

Text as Data

Justin Grimmer

Associate Professor
Department of Political Science
Stanford University

September 25th, 2014

A General Framework for the Class

A focus on **tasks**

A General Framework for the Class

A focus on **tasks**

- 1) Decide on some objective to accomplish

A General Framework for the Class

A focus on **tasks**

- 1) Decide on some objective to accomplish
 - Classify

A General Framework for the Class

A focus on **tasks**

- 1) Decide on some objective to accomplish
 - Classify
 - Cluster

A General Framework for the Class

A focus on **tasks**

1) Decide on some objective to accomplish

- Classify
- Cluster
- Predict

A General Framework for the Class

A focus on **tasks**

1) Decide on some objective to accomplish

- Classify
- Cluster
- Predict
- Describe

A General Framework for the Class

A focus on **tasks**

1) Decide on some objective to accomplish

- Classify
- Cluster
- Predict
- Describe
- Measure covariance, discover latent structure, find nearest neighbor, ...

A General Framework for the Class

- 2) Use an objective function \rightsquigarrow measure performance at task

A General Framework for the Class

- 2) Use an objective function \rightsquigarrow measure performance at task
Suppose we have data $\mathbf{X} \in \mathcal{X}$ and parameters $\boldsymbol{\theta} \in \Theta$

A General Framework for the Class

- 2) Use an objective function \rightsquigarrow measure performance at task
Suppose we have data $\mathbf{X} \in \mathcal{X}$ and parameters $\boldsymbol{\theta} \in \Theta$

$$f : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$$

A General Framework for the Class

- 2) Use an objective function \rightsquigarrow measure performance at task
Suppose we have data $\mathbf{X} \in \mathcal{X}$ and parameters $\boldsymbol{\theta} \in \Theta$

$$f : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$$

$$f(\mathbf{X}, \boldsymbol{\theta})$$

A General Framework for the Class

- 2) Use an objective function \rightsquigarrow measure performance at task
Suppose we have data $\mathbf{X} \in \mathcal{X}$ and parameters $\boldsymbol{\theta} \in \Theta$

$$f : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$$
$$f(\mathbf{X}, \boldsymbol{\theta})$$

- Derived from task/used for statistical properties \rightsquigarrow minimize sum of squared residuals

A General Framework for the Class

- 2) Use an objective function \rightsquigarrow measure performance at task
Suppose we have data $\mathbf{X} \in \mathcal{X}$ and parameters $\boldsymbol{\theta} \in \Theta$

$$f : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$$
$$f(\mathbf{X}, \boldsymbol{\theta})$$

- Derived from task/used for statistical properties \rightsquigarrow minimize sum of squared residuals
- Formalization of intuition about “good” performance \rightsquigarrow k-means clustering

A General Framework for the Class

- 2) Use an objective function \rightsquigarrow measure performance at task
Suppose we have data $\mathbf{X} \in \mathcal{X}$ and parameters $\boldsymbol{\theta} \in \Theta$

$$f : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$$
$$f(\mathbf{X}, \boldsymbol{\theta})$$

- Derived from task/used for statistical properties \rightsquigarrow minimize sum of squared residuals
- Formalization of intuition about “good” performance \rightsquigarrow k-means clustering
- Data generating process

A general Framework for the Class

3) Method for optimizing objective function.

A general Framework for the Class

- 3) Method for optimizing objective function.
Find $\theta^* \in \Theta$ such that

A general Framework for the Class

- 3) Method for optimizing objective function.
Find $\theta^* \in \Theta$ such that

$$f(\mathbf{X}, \theta^*)$$

A general Framework for the Class

- 3) Method for optimizing objective function.
Find $\theta^* \in \Theta$ such that

$$f(\mathbf{X}, \theta^*)$$

is a maximum (minimum)

A general Framework for the Class

- 3) Method for optimizing objective function.
Find $\theta^* \in \Theta$ such that

$$f(\mathbf{X}, \theta^*)$$

is a maximum (minimum)

- Analytic methods (Calculus)

A general Framework for the Class

- 3) Method for optimizing objective function.
Find $\theta^* \in \Theta$ such that

$$f(\mathbf{X}, \theta^*)$$

is a maximum (minimum)

- Analytic methods (Calculus)
- **Computational** methods

Today: Using (Bayesian) Statistics to Obtain Objective Functions

- Encode assumptions in **data generating process** \rightsquigarrow hierarchical model
- Assume parameters and data are random variables
- Conditional probability statement \rightsquigarrow objective function
- Use computational tools to optimize objective function

Today: Using (Bayesian) Statistics to Obtain Objective Functions

Plan of Attack:

- 1) Write out any joint density function as conditional relationship
- 2) Show how this can be an objective function **even if you've never taken likelihood/Bayesian/...**
- 3) Discuss how to **computationally** optimize

Joint Distributions of Random Variables

Definition

Suppose that we have random variables $\mathbf{X} = (X_1, X_2, \dots, X_K)$. We will say that \mathbf{X} is a jointly continuous random variable if for all $\mathbf{X} \in \mathbb{R}^K$ there exists a function $f : \mathbb{R}^K \rightarrow \mathbb{R}$ such that for all $C \subset \mathbb{R}^K$,

$$P(\mathbf{X} \in C) = \iint \dots \iint_{(\mathbf{x}) \in C} f(\mathbf{x}) d\mathbf{X}$$

- A joint density $f(\mathbf{x}) = f(x_1, x_2, \dots, x_K)$ encodes information about the behavior of the random variable \mathbf{X}

Marginal Distribution

Definition

Suppose $\mathbf{X} = (X_1, X_2, \dots, X_K)$ is a jointly continuous random variable. Define $f_{X_j}(x)$ as the **marginal** probability density function for X_j ,

$$\begin{aligned} f_{X_j}(x) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(\mathbf{x}) dx_1 dx_2 \dots dx_{j-1} dx_{j+1} \dots dx_{K-1} dx_K \\ &= \int_{\mathbb{R}^{K-1}} f(\mathbf{x}) d\mathbf{X}_{-j} \end{aligned}$$

- To obtain the marginal distribution, $f_{X_j}(x)$ we integrate over all dimensions but j

Conditional Distributions and Independence of Random Variables

Definition

Suppose \mathbf{X} is a jointly continuous random variable. Define $f_{\mathbf{X}_{-j}|X_j}(\mathbf{x})$ as the conditional density function,

$$f_{\mathbf{X}_{-j}|X_j}(\mathbf{x}_{-j}|x_j) = \frac{f(x_1, x_2, \dots, x_K)}{f_{X_j}(x_j)}$$

Conditional Distributions and Independence of Random Variables

Definition

Suppose \mathbf{X} is a jointly continuous random variable. Define $f_{\mathbf{X}_{-j}|\mathbf{X}_j}(\mathbf{x})$ as the conditional density function,

$$f_{\mathbf{X}_{-j}|\mathbf{X}_j}(\mathbf{x}_{-j}|\mathbf{x}_j) = \frac{f(x_1, x_2, \dots, x_K)}{f_{\mathbf{X}_j}(\mathbf{x}_j)}$$

Two random variables X_1 and X_2 are independent if

$$f_{X_1, X_2}(x_1, x_2) = f_{X_1}(x_1)f_{X_2}(x_2)$$

Conditional Independence of Random Variables

Definition

Two random variables X_1 and X_2 are conditionally independent given X_3 if

$$f_{(X_1, X_2)|X_3}(x_1, x_2|x_3) = f_{X_1|X_3}(x_1|x_3)f_{X_2|X_3}(x_2|x_3)$$

Conditional Independence of Random Variables

Definition

Two random variables X_1 and X_2 are conditionally independent given X_3 if

$$f_{(X_1, X_2)|X_3}(x_1, x_2|x_3) = f_{X_1|X_3}(x_1|x_3)f_{X_2|X_3}(x_2|x_3)$$

Equivalently,

$$f_{X_1|X_2, X_3}(x_1|x_2, x_3) = f_{X_1|X_3}(x_1|x_3)$$

Joint Density as Conditional Relationship

Theorem

Suppose $\mathbf{X} = (X_1, X_2, \dots, X_K)$ is a jointly continuous random variable. Then

$$\begin{aligned} f(\mathbf{x}) &= f(x_1, x_2, \dots, x_K) \\ &= f_{X_1}(x_1) f_{X_2|X_1}(x_2|x_1) f_{X_3|X_2X_1}(x_3|x_2, x_1) \dots f_{X_K|X_1 \dots X_{K-1}}(x_K|x_1, \dots, x_{K-1}) \end{aligned}$$

Joint Density as Conditional Relationship

Theorem

Suppose $\mathbf{X} = (X_1, X_2, \dots, X_K)$ is a jointly continuous random variable. Then

$$\begin{aligned} f(\mathbf{x}) &= f(x_1, x_2, \dots, x_K) \\ &= f_{X_1}(x_1) f_{X_2|X_1}(x_2|x_1) f_{X_3|X_2X_1}(x_3|x_2, x_1) \dots f_{X_K|X_1 \dots X_{K-1}}(x_K|x_1, \dots, x_{K-1}) \end{aligned}$$

- We can write joint distributions as a product of conditional distributions

Joint Density as Conditional Relationship

Theorem

Suppose $\mathbf{X} = (X_1, X_2, \dots, X_K)$ is a jointly continuous random variable. Then

$$\begin{aligned} f(\mathbf{x}) &= f(x_1, x_2, \dots, x_K) \\ &= f_{X_1}(x_1) f_{X_2|X_1}(x_2|x_1) f_{X_3|X_2X_1}(x_3|x_2, x_1) \dots f_{X_K|X_1 \dots X_{K-1}}(x_K|x_1, \dots, x_{K-1}) \end{aligned}$$

- We can write joint distributions as a product of conditional distributions
- If there are **conditional independences** in density we can simplify \rightsquigarrow simplify some conditional expressions

Beta-Binomial Model

Definition

Suppose Y is a continuous random variable with $Y \in [0, 1]$ and pdf of Y given by

$$f(y) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} y^{\alpha_1-1} (1-y)^{\alpha_2-1}$$

*Then we will say Y is a **Beta** distribution with parameters α_1 and α_2 . Equivalently,*

$$Y \sim \text{Beta}(\alpha_1, \alpha_2)$$

Beta-Binomial Model

Definition

Suppose Y is a continuous random variable with $Y \in [0, 1]$ and pdf of Y given by

$$f(y) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} y^{\alpha_1-1} (1-y)^{\alpha_2-1}$$

*Then we will say Y is a **Beta** distribution with parameters α_1 and α_2 . Equivalently,*

$$Y \sim \text{Beta}(\alpha_1, \alpha_2)$$

- Beta is a distribution on **proportions**

Beta-Binomial Model

Definition

Suppose Y is a continuous random variable with $Y \in [0, 1]$ and pdf of Y given by

$$f(y) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} y^{\alpha_1-1} (1-y)^{\alpha_2-1}$$

*Then we will say Y is a **Beta** distribution with parameters α_1 and α_2 . Equivalently,*

$$Y \sim \text{Beta}(\alpha_1, \alpha_2)$$

- Beta is a distribution on **proportions**
- Beta is a special case of the **Dirichlet** distribution

Beta-Binomial Model

Definition

Suppose Y is a continuous random variable with $Y \in [0, 1]$ and pdf of Y given by

$$f(y) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} y^{\alpha_1-1} (1-y)^{\alpha_2-1}$$

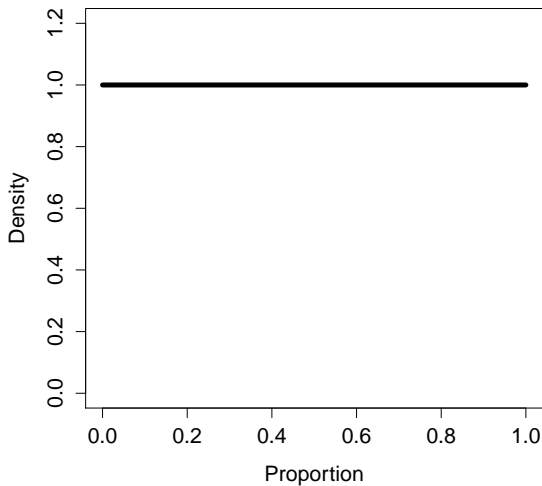
*Then we will say Y is a **Beta** distribution with parameters α_1 and α_2 . Equivalently,*

$$Y \sim \text{Beta}(\alpha_1, \alpha_2)$$

- Beta is a distribution on **proportions**
- Beta is a special case of the **Dirichlet** distribution
- $E[Y] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$

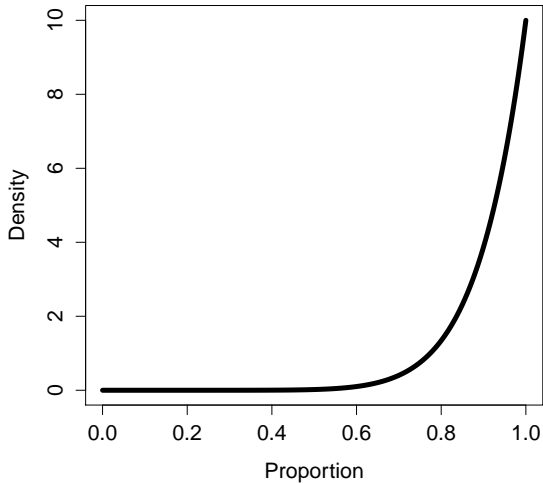
Beta Distribution

$\alpha_1 = 1, \alpha_2 = 1$



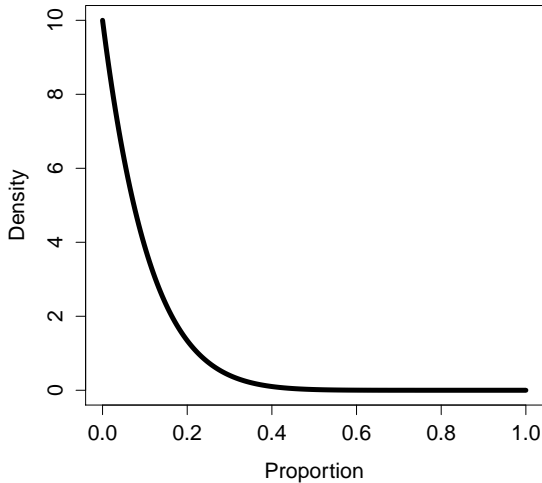
Beta Distribution

$\alpha_1 = 10$, $\alpha_2 = 1$



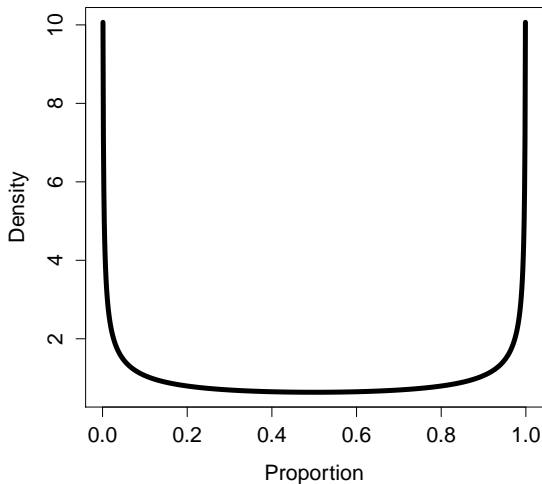
Beta Distribution

$\alpha_1 = 1, \alpha_2 = 10$



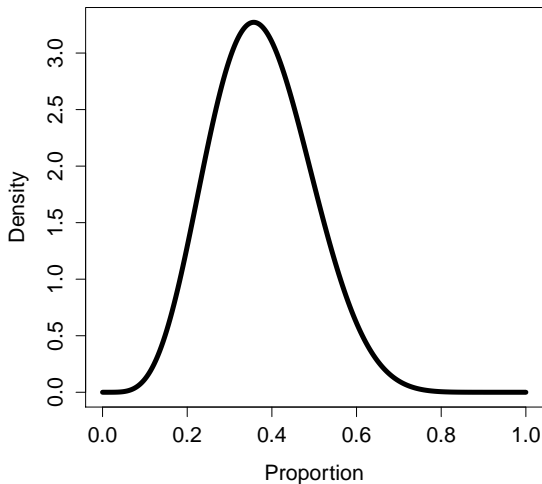
Beta Distribution

$\alpha_1 = 0.5, \alpha_2 = 0.5$



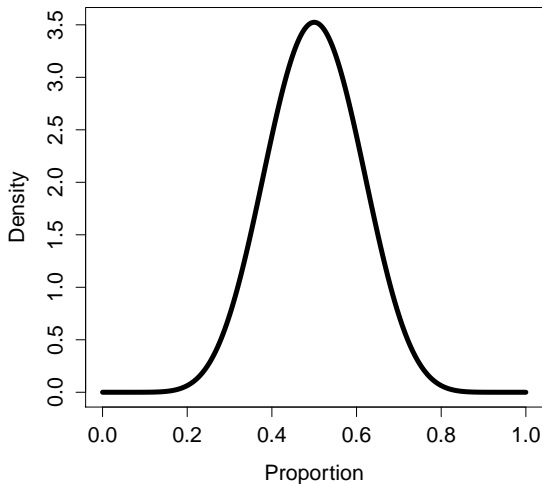
Beta Distribution

$\alpha_1 = 6, \alpha_2 = 10$



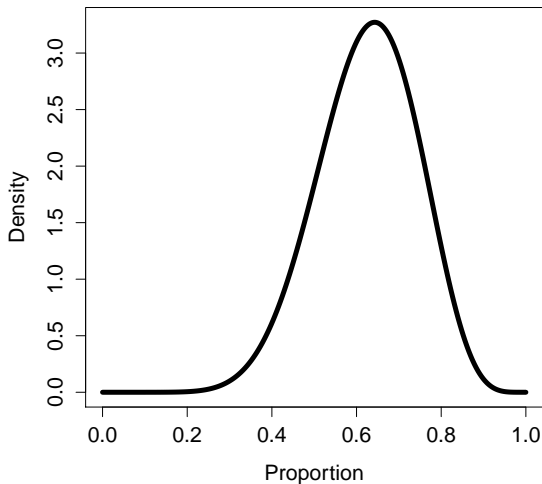
Beta Distribution

$\alpha_1 = 10, \alpha_2 = 10$



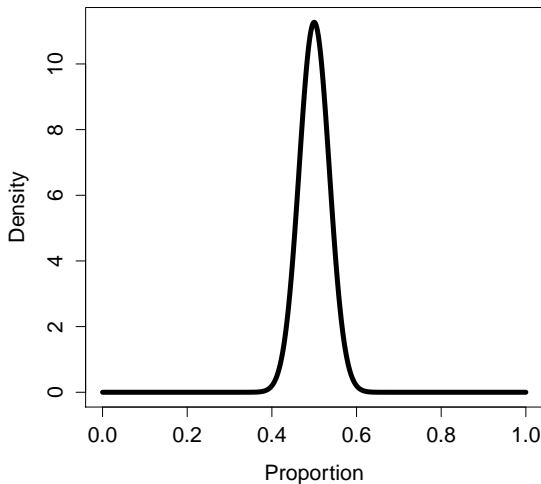
Beta Distribution

alpha_1 = 10, alpha_2 = 6



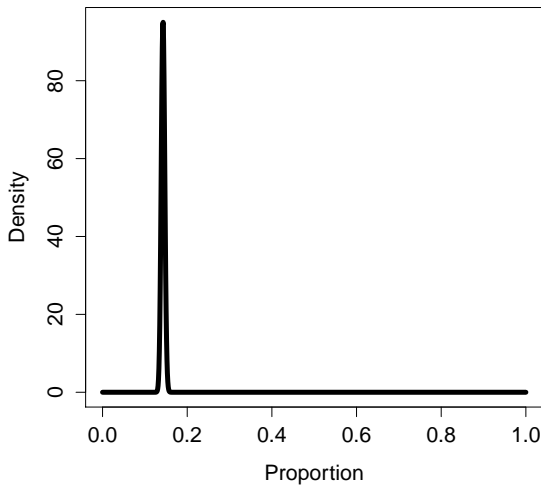
Beta Distribution

$\alpha_1 = 100, \alpha_2 = 100$



Beta Distribution

$\alpha_1 = 1000$, $\alpha_2 = 6000$



Beta-Binomial Model

Suppose we are interested in the use of the word “Obamacare” in a legislator’s statements

Beta-Binomial Model

Suppose we are interested in the use of the word “Obamacare” in a legislator’s statements

We want to infer π , we observe incidence in a statement

Beta-Binomial Model

Suppose we are interested in the use of the word “Obamacare” in a legislator’s statements

We want to infer π , we observe incidence in a statement

$$\pi \sim \text{Beta}(\alpha_1, \alpha_2)$$

Beta-Binomial Model

Suppose we are interested in the use of the word “Obamacare” in a legislator’s statements

We want to infer π , we observe incidence in a statement

$$\begin{aligned}\pi &\sim \text{Beta}(\alpha_1, \alpha_2) \\ Y|\pi &\sim \text{Bernoulli}(\pi)\end{aligned}$$

Beta-Binomial Model

Suppose we are interested in the use of the word “Obamacare” in a legislator’s statements

We want to infer π , we observe incidence in a statement

$$\begin{aligned}\pi &\sim \text{Beta}(\alpha_1, \alpha_2) \\ Y|\pi &\sim \text{Bernoulli}(\pi)\end{aligned}$$

We can write $p(\pi, y|\alpha_1, \alpha_2)$ as

Beta-Binomial Model

Suppose we are interested in the use of the word “Obamacare” in a legislator’s statements

We want to infer π , we observe incidence in a statement

$$\begin{aligned}\pi &\sim \text{Beta}(\alpha_1, \alpha_2) \\ Y|\pi &\sim \text{Bernoulli}(\pi)\end{aligned}$$

We can write $p(\pi, y|\alpha_1, \alpha_2)$ as

$$p(\pi, y|\alpha_1, \alpha_2)$$

Beta-Binomial Model

Suppose we are interested in the use of the word “Obamacare” in a legislator’s statements

We want to infer π , we observe incidence in a statement

$$\begin{aligned}\pi &\sim \text{Beta}(\alpha_1, \alpha_2) \\ Y|\pi &\sim \text{Bernoulli}(\pi)\end{aligned}$$

We can write $p(\pi, y|\alpha_1, \alpha_2)$ as

$$p(\pi, y|\alpha_1, \alpha_2) = p(\pi|\alpha_1, \alpha_2)p(y|\pi, \alpha_1, \alpha_2) = p(\pi|\alpha_1, \alpha_2)p(y|\pi)$$

Beta-Binomial Model

Suppose we are interested in the use of the word “Obamacare” in a legislator’s statements

We want to infer π , we observe incidence in a statement

$$\begin{aligned}\pi &\sim \text{Beta}(\alpha_1, \alpha_2) \\ Y|\pi &\sim \text{Bernoulli}(\pi)\end{aligned}$$

We can write $p(\pi, y|\alpha_1, \alpha_2)$ as

$$\begin{aligned}p(\pi, y|\alpha_1, \alpha_2) &= p(\pi|\alpha_1, \alpha_2)p(y|\pi, \alpha_1, \alpha_2) = p(\pi|\alpha_1, \alpha_2)p(y|\pi) \\ &= \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)}\pi^{\alpha_1-1}(1 - \pi)^{\alpha_2-1}\pi^y(1 - \pi)^{1-y}\end{aligned}$$

Beta-Binomial Model

Suppose we are interested in the use of the word “Obamacare” in a legislator’s statements

We want to infer π , we observe incidence in a statement

$$\begin{aligned}\pi &\sim \text{Beta}(\alpha_1, \alpha_2) \\ Y|\pi &\sim \text{Bernoulli}(\pi)\end{aligned}$$

We can write $p(\pi, y|\alpha_1, \alpha_2)$ as

$$\begin{aligned}p(\pi, y|\alpha_1, \alpha_2) &= p(\pi|\alpha_1, \alpha_2)p(y|\pi, \alpha_1, \alpha_2) = p(\pi|\alpha_1, \alpha_2)p(y|\pi) \\ &= \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \pi^{\alpha_1-1} (1 - \pi)^{\alpha_2-1} \pi^y (1 - \pi)^{1-y} \\ &= \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \pi^{\alpha_1+y-1} (1 - \pi)^{\alpha_2+1-y-1}\end{aligned}$$

Beta-Binomial Model

Suppose now we observe n statements

$$\pi \sim \text{Beta}(\alpha_1, \alpha_2)$$

$$Y_i | \pi \sim \text{Bernoulli}(\pi) \text{ for } i = 1, \dots, n$$

$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$$

Beta-Binomial Model

Suppose now we observe n statements

$$\pi \sim \text{Beta}(\alpha_1, \alpha_2)$$

$$Y_i | \pi \sim \text{Bernoulli}(\pi) \text{ for } i = 1, \dots, n$$

$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$$

We can write $p(\pi, \mathbf{y} | \alpha_1, \alpha_2)$ as

Beta-Binomial Model

Suppose now we observe n statements

$$\pi \sim \text{Beta}(\alpha_1, \alpha_2)$$

$$Y_i | \pi \sim \text{Bernoulli}(\pi) \text{ for } i = 1, \dots, n$$

$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$$

We can write $p(\pi, \mathbf{y} | \alpha_1, \alpha_2)$ as

$$p(\pi, \mathbf{y} | \alpha_1, \alpha_2)$$

Beta-Binomial Model

Suppose now we observe n statements

$$\pi \sim \text{Beta}(\alpha_1, \alpha_2)$$

$$Y_i | \pi \sim \text{Bernoulli}(\pi) \text{ for } i = 1, \dots, n$$

$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$$

We can write $p(\pi, \mathbf{y} | \alpha_1, \alpha_2)$ as

$$p(\pi, \mathbf{y} | \alpha_1, \alpha_2) = p(\pi | \alpha_1, \alpha_2) p(\mathbf{y} | \pi)$$

Beta-Binomial Model

Suppose now we observe n statements

$$\pi \sim \text{Beta}(\alpha_1, \alpha_2)$$

$$Y_i | \pi \sim \text{Bernoulli}(\pi) \text{ for } i = 1, \dots, n$$

$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$$

We can write $p(\pi, \mathbf{y} | \alpha_1, \alpha_2)$ as

$$\begin{aligned} p(\pi, \mathbf{y} | \alpha_1, \alpha_2) &= p(\pi | \alpha_1, \alpha_2) p(\mathbf{y} | \pi) \\ &= p(\pi | \alpha_1, \alpha_2) \prod_{i=1}^n p(y_i | \pi) \end{aligned}$$

Beta-Binomial Model

Suppose now we observe n statements

$$\begin{aligned}\pi &\sim \text{Beta}(\alpha_1, \alpha_2) \\ Y_i | \pi &\sim \text{Bernoulli}(\pi) \text{ for } i = 1, \dots, n\end{aligned}$$

$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$$

We can write $p(\pi, \mathbf{y} | \alpha_1, \alpha_2)$ as

$$\begin{aligned}p(\pi, \mathbf{y} | \alpha_1, \alpha_2) &= p(\pi | \alpha_1, \alpha_2) p(\mathbf{y} | \pi) \\ &= p(\pi | \alpha_1, \alpha_2) \prod_{i=1}^n p(y_i | \pi) \\ &= \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1) \Gamma(\alpha_2)} \pi^{\alpha_1 - 1} (1 - \pi)^{\alpha_2 - 1} \prod_{i=1}^n \pi^{y_i} (1 - \pi)^{1 - y_i}\end{aligned}$$

Beta-Binomial Model

Suppose now we observe n statements

$$\begin{aligned}\pi &\sim \text{Beta}(\alpha_1, \alpha_2) \\ Y_i | \pi &\sim \text{Bernoulli}(\pi) \text{ for } i = 1, \dots, n\end{aligned}$$

$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$$

We can write $p(\pi, \mathbf{y} | \alpha_1, \alpha_2)$ as

$$\begin{aligned}p(\pi, \mathbf{y} | \alpha_1, \alpha_2) &= p(\pi | \alpha_1, \alpha_2) p(\mathbf{y} | \pi) \\ &= p(\pi | \alpha_1, \alpha_2) \prod_{i=1}^n p(y_i | \pi) \\ &= \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1) \Gamma(\alpha_2)} \pi^{\alpha_1 - 1} (1 - \pi)^{\alpha_2 - 1} \prod_{i=1}^n \pi^{y_i} (1 - \pi)^{1 - y_i} \\ &= \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1) \Gamma(\alpha_2)} \pi^{\sum_{i=1}^n y_i + \alpha_1 - 1} (1 - \pi)^{n + \alpha_2 - \sum_{i=1}^n y_i - 1}\end{aligned}$$

Bayes' Theorem for Continuous Random Variables

Theorem

Suppose we have jointly continuous random variables X_1 and X_2 . Then,

$$f_{X_1|X_2}(x_1|x_2) = \frac{f_{X_1}(x_1)f_{X_2|X_1}(x_2|x_1)}{f_{X_2}(x_2)}$$

Beta-Binomial Model

We observe \mathbf{y} and we want to learn about π

Beta-Binomial Model

We observe \mathbf{y} and we want to learn about π

Condition on data, describe function of π .

Beta-Binomial Model

We observe \mathbf{y} and we want to learn about π

Condition on data, describe function of π .

$$p(\pi|\mathbf{y}, \alpha_1, \alpha_2) = \frac{p(\pi|\alpha_1, \alpha_2)p(\mathbf{y}|\pi, \alpha_1, \alpha_2)}{p(\mathbf{y})}$$

Beta-Binomial Model

We observe \mathbf{y} and we want to learn about π

Condition on data, describe function of π .

$$\begin{aligned} p(\pi|\mathbf{y}, \alpha_1, \alpha_2) &= \frac{p(\pi|\alpha_1, \alpha_2)p(\mathbf{y}|\pi, \alpha_1, \alpha_2)}{p(\mathbf{y})} \\ &\propto p(\pi|\alpha_1, \alpha_2)p(\mathbf{y}|\pi) \end{aligned}$$

Beta-Binomial Model

We observe \mathbf{y} and we want to learn about π

Condition on data, describe function of π .

$$\begin{aligned} p(\pi|\mathbf{y}, \alpha_1, \alpha_2) &= \frac{p(\pi|\alpha_1, \alpha_2)p(\mathbf{y}|\pi, \alpha_1, \alpha_2)}{p(\mathbf{y})} \\ &\propto p(\pi|\alpha_1, \alpha_2)p(\mathbf{y}|\pi) \\ &\propto \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \pi^{\sum_{i=1}^n y_i + \alpha_1 - 1} (1 - \pi)^{n + \alpha_2 - \sum_{i=1}^n y_i - 1} \end{aligned}$$

Beta-Binomial Model

We observe \mathbf{y} and we want to learn about π

Condition on data, describe function of π .

$$\begin{aligned} p(\pi|\mathbf{y}, \alpha_1, \alpha_2) &= \frac{p(\pi|\alpha_1, \alpha_2)p(\mathbf{y}|\pi, \alpha_1, \alpha_2)}{p(\mathbf{y})} \\ &\propto p(\pi|\alpha_1, \alpha_2)p(\mathbf{y}|\pi) \\ &\propto \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \pi^{\sum_{i=1}^n y_i + \alpha_1 - 1} (1 - \pi)^{n + \alpha_2 - \sum_{i=1}^n y_i - 1} \end{aligned}$$

Defines a function of π , which we can use to describe the data.

Beta-Binomial Model

We observe \mathbf{y} and we want to learn about π

Condition on data, describe function of π .

$$\begin{aligned} p(\pi|\mathbf{y}, \alpha_1, \alpha_2) &= \frac{p(\pi|\alpha_1, \alpha_2)p(\mathbf{y}|\pi, \alpha_1, \alpha_2)}{p(\mathbf{y})} \\ &\propto p(\pi|\alpha_1, \alpha_2)p(\mathbf{y}|\pi) \\ &\propto \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \pi^{\sum_{i=1}^n y_i + \alpha_1 - 1} (1 - \pi)^{n + \alpha_2 - \sum_{i=1}^n y_i - 1} \end{aligned}$$

Defines a function of π , which we can use to describe the data.

Optimize \rightsquigarrow **analytically** or computationally.

Theorem

Suppose $f : \mathbb{R}^K \rightarrow \mathbb{R}_+$. If \mathbf{x}^ maximizes $\log(f(\mathbf{x}))$, then \mathbf{x}^* maximizes $f(\mathbf{x})$.*

Theorem

Suppose $f : \mathbb{R}^K \rightarrow \mathbb{R}_+$. If \mathbf{x}^* maximizes $\log(f(\mathbf{x}))$, then \mathbf{x}^* maximizes $f(\mathbf{x})$.

$$\log(p(\pi|\mathbf{y}, \alpha_1, \alpha_2))$$

Theorem

Suppose $f : \mathbb{R}^K \rightarrow \mathbb{R}_+$. If \mathbf{x}^* maximizes $\log(f(\mathbf{x}))$, then \mathbf{x}^* maximizes $f(\mathbf{x})$.

$$\begin{aligned} \log(p(\pi|\mathbf{y}, \alpha_1, \alpha_2)) &= \left(\sum_{i=1}^n y_i + \alpha_1 - 1 \right) \log \pi \\ &+ (n + \alpha_2 - \sum_{i=1}^n y_i - 1)(1 - \log \pi) + \textcolor{red}{c} \end{aligned}$$

Theorem

Suppose $f : \mathbb{R}^K \rightarrow \mathbb{R}_+$. If \mathbf{x}^* maximizes $\log(f(\mathbf{x}))$, then \mathbf{x}^* maximizes $f(\mathbf{x})$.

$$\log(p(\pi|\mathbf{y}, \alpha_1, \alpha_2)) = \left(\sum_{i=1}^n y_i + \alpha_1 - 1 \right) \log \pi$$

$$+ (n + \alpha_2 - \sum_{i=1}^n y_i - 1)(1 - \log \pi) + \textcolor{red}{c}$$

$$\frac{\partial \log(p(\pi|\mathbf{y}, \alpha_1, \alpha_2))}{\partial \pi} = \frac{\sum_{i=1}^n y_i + \alpha_1 - 1}{\pi} - \frac{(n + \alpha_2 - \sum_{i=1}^n y_i - 1)}{1 - \pi}$$

Theorem

Suppose $f : \mathbb{R}^K \rightarrow \mathbb{R}_+$. If \mathbf{x}^* maximizes $\log(f(\mathbf{x}))$, then \mathbf{x}^* maximizes $f(\mathbf{x})$.

$$\log(p(\pi|\mathbf{y}, \alpha_1, \alpha_2)) = \left(\sum_{i=1}^n y_i + \alpha_1 - 1 \right) \log \pi$$

$$+ (n + \alpha_2 - \sum_{i=1}^n y_i - 1)(1 - \log \pi) + c$$

$$\frac{\partial \log(p(\pi|\mathbf{y}, \alpha_1, \alpha_2))}{\partial \pi} = \frac{\sum_{i=1}^n y_i + \alpha_1 - 1}{\pi} - \frac{(n + \alpha_2 - \sum_{i=1}^n y_i - 1)}{1 - \pi}$$

$$0 = \frac{\sum_{i=1}^n y_i + \alpha_1 - 1}{\pi^*} - \frac{n + \alpha_2 - \sum_{i=1}^n y_i - 1}{1 - \pi^*}$$

Theorem

Suppose $f : \mathbb{R}^K \rightarrow \mathbb{R}_+$. If \mathbf{x}^* maximizes $\log(f(\mathbf{x}))$, then \mathbf{x}^* maximizes $f(\mathbf{x})$.

$$\log(p(\pi|\mathbf{y}, \alpha_1, \alpha_2)) = \left(\sum_{i=1}^n y_i + \alpha_1 - 1 \right) \log \pi$$

$$+ (n + \alpha_2 - \sum_{i=1}^n y_i - 1)(1 - \log \pi) + c$$

$$\frac{\partial \log(p(\pi|\mathbf{y}, \alpha_1, \alpha_2))}{\partial \pi} = \frac{\sum_{i=1}^n y_i + \alpha_1 - 1}{\pi} - \frac{(n + \alpha_2 - \sum_{i=1}^n y_i - 1)}{1 - \pi}$$

$$0 = \frac{\sum_{i=1}^n y_i + \alpha_1 - 1}{\pi^*} - \frac{n + \alpha_2 - \sum_{i=1}^n y_i - 1}{1 - \pi^*}$$

$$\pi^* = \frac{\sum_{i=1}^n y_i + \alpha_1 - 1}{n + \alpha_1 + \alpha_2 - 2}$$

Theorem

Suppose $f : \mathbb{R}^K \rightarrow \mathbb{R}_+$. If \mathbf{x}^* maximizes $\log(f(\mathbf{x}))$, then \mathbf{x}^* maximizes $f(\mathbf{x})$.

$$\begin{aligned}\log(p(\pi|\mathbf{y}, \alpha_1, \alpha_2)) &= \left(\sum_{i=1}^n y_i + \alpha_1 - 1 \right) \log \pi \\ &\quad + (n + \alpha_2 - \sum_{i=1}^n y_i - 1)(1 - \log \pi) + c \\ \frac{\partial \log(p(\pi|\mathbf{y}, \alpha_1, \alpha_2))}{\partial \pi} &= \frac{\sum_{i=1}^n y_i + \alpha_1 - 1}{\pi} - \frac{(n + \alpha_2 - \sum_{i=1}^n y_i - 1)}{1 - \pi} \\ 0 &= \frac{\sum_{i=1}^n y_i + \alpha_1 - 1}{\pi^*} - \frac{n + \alpha_2 - \sum_{i=1}^n y_i - 1}{1 - \pi^*} \\ \pi^* &= \frac{\sum_{i=1}^n y_i + \alpha_1 - 1}{n + \alpha_1 + \alpha_2 - 2}\end{aligned}$$

Second derivative test \rightsquigarrow maximum

The Probit Model

Suppose if we're interested in regression

The Probit Model

Suppose if we're interested in regression \rightsquigarrow prediction, classification, description

The Probit Model

Suppose if we're interested in regression \rightsquigarrow prediction, classification, description

Assume the following data generation process

The Probit Model

Suppose if we're interested in regression \rightsquigarrow prediction, classification, description

Assume the following data generation process

$$Y_i \sim \text{Bernoulli}(\pi_i)$$

The Probit Model

Suppose if we're interested in regression \rightsquigarrow prediction, classification, description

Assume the following data generation process

$$Y_i \sim \text{Bernoulli}(\pi_i)$$

$$\pi_i = \Phi(\mathbf{X}_i\beta)$$

The Probit Model

Suppose if we're interested in regression \rightsquigarrow prediction, classification, description

Assume the following data generation process

$$Y_i \sim \text{Bernoulli}(\pi_i)$$

$$\pi_i = \Phi(\mathbf{X}_i\beta)$$

where

The Probit Model

Suppose if we're interested in regression \rightsquigarrow prediction, classification, description

Assume the following data generation process

$$\begin{aligned} Y_i &\sim \text{Bernoulli}(\pi_i) \\ \pi_i &= \Phi(\mathbf{X}_i\beta) \end{aligned}$$

where

$\Phi(\cdot)$ is the cumulative normal distribution function

The Probit Model

Suppose if we're interested in regression \rightsquigarrow prediction, classification, description

Assume the following data generation process

$$\begin{aligned} Y_i &\sim \text{Bernoulli}(\pi_i) \\ \pi_i &= \Phi(\mathbf{X}_i\beta) \end{aligned}$$

where

$\Phi(\cdot)$ is the cumulative normal distribution function

$$\mathbf{X}_i = (1, x_i)$$

The Probit Model

Suppose if we're interested in regression \rightsquigarrow prediction, classification, description

Assume the following data generation process

$$\begin{aligned} Y_i &\sim \text{Bernoulli}(\pi_i) \\ \pi_i &= \Phi(\mathbf{X}_i\boldsymbol{\beta}) \end{aligned}$$

where

$\Phi(\cdot)$ is the cumulative normal distribution function

$$\mathbf{X}_i = (1, x_i)$$

$$\boldsymbol{\beta} = (\beta_1, \beta_2) \rightsquigarrow \text{parameters.}$$

The Probit Model

Suppose if we're interested in regression \rightsquigarrow prediction, classification, description

Assume the following data generation process

$$\begin{aligned} Y_i &\sim \text{Bernoulli}(\pi_i) \\ \pi_i &= \Phi(\mathbf{X}_i\boldsymbol{\beta}) \end{aligned}$$

where

$\Phi(\cdot)$ is the cumulative normal distribution function

$$\mathbf{X}_i = (1, x_i)$$

$$\boldsymbol{\beta} = (\beta_1, \beta_2) \rightsquigarrow \text{parameters.}$$

N observations

The Probit Model

Suppose if we're interested in regression \rightsquigarrow prediction, classification, description

Assume the following data generation process

$$\begin{aligned}Y_i &\sim \text{Bernoulli}(\pi_i) \\ \pi_i &= \Phi(\mathbf{X}_i\boldsymbol{\beta})\end{aligned}$$

where

$\Phi(\cdot)$ is the cumulative normal distribution function

$$\mathbf{X}_i = (1, x_i)$$

$$\boldsymbol{\beta} = (\beta_1, \beta_2) \rightsquigarrow \text{parameters.}$$

N observations

Implicit (improper) prior

The Probit Model \rightsquigarrow Objective Function

$$\begin{aligned} p(\beta|\mathbf{y}, \mathbf{X}) &\propto p(\mathbf{y}|\beta, \mathbf{X}) \\ &\propto \prod_{i=1}^N \Phi(\mathbf{X}_i\beta)^{y_i} (1 - \Phi(\mathbf{X}_i\beta))^{1-y_i} \end{aligned}$$

The Probit Model \rightsquigarrow Optimization

Theorem

Suppose $f : \mathbb{R}^K \rightarrow \mathbb{R}_+$. If \mathbf{x}^ maximizes $\log(f(\mathbf{x}))$, then \mathbf{x}^* maximizes $f(\mathbf{x})$.*

$$\log(p(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X})) = \sum_{i=1}^N [y_i \log(\Phi(\mathbf{X}_i\boldsymbol{\beta})) + (1 - y_i) \log(1 - \Phi(\mathbf{X}_i\boldsymbol{\beta}))] + \textcolor{red}{c}$$

Find $\boldsymbol{\beta}^*$ to maximize $\log(p(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X})) \rightsquigarrow$ computational method

Computational Optimization Approaches

Analytic (Closed form) \rightsquigarrow Often difficult, impractical, or unavailable

Computational Optimization Approaches

Analytic (Closed form) \rightsquigarrow Often difficult, impractical, or unavailable
Computational

Computational Optimization Approaches

Analytic (Closed form) \rightsquigarrow Often difficult, impractical, or unavailable
Computational \rightsquigarrow **iterative** algorithm that converges to a solution
(hopefully the right one!)

Computational Optimization Approaches

- Analytic** (Closed form) \rightsquigarrow Often difficult, impractical, or unavailable
- Computational** \rightsquigarrow **iterative** algorithm that converges to a solution (hopefully the right one!)
- Methods for optimization:

Computational Optimization Approaches

Analytic (Closed form) \rightsquigarrow Often difficult, impractical, or unavailable
Computational \rightsquigarrow **iterative** algorithm that converges to a solution
(hopefully the right one!)

- Methods for optimization:
 - **Newton's method** and **BFGS**

Computational Optimization Approaches

Analytic (Closed form) \rightsquigarrow Often difficult, impractical, or unavailable
Computational \rightsquigarrow **iterative** algorithm that converges to a solution
(hopefully the right one!)

- Methods for optimization:
 - **Newton's method** and **BFGS**
 - Gradient descent (ascent)

Computational Optimization Approaches

Analytic (Closed form) \rightsquigarrow Often difficult, impractical, or unavailable
Computational \rightsquigarrow **iterative** algorithm that converges to a solution
(hopefully the right one!)

- Methods for optimization:
 - **Newton's method** and **BFGS**
 - Gradient descent (ascent)
 - Expectation Maximization

Computational Optimization Approaches

Analytic (Closed form) \rightsquigarrow Often difficult, impractical, or unavailable
Computational \rightsquigarrow **iterative** algorithm that converges to a solution
(hopefully the right one!)

- Methods for optimization:
 - **Newton's method** and **BFGS**
 - Gradient descent (ascent)
 - Expectation Maximization
 - Genetic Optimization

Computational Optimization Approaches

Analytic (Closed form) \rightsquigarrow Often difficult, impractical, or unavailable
Computational \rightsquigarrow **iterative** algorithm that converges to a solution
(hopefully the right one!)

- Methods for optimization:
 - **Newton's method** and **BFGS**
 - Gradient descent (ascent)
 - Expectation Maximization
 - Genetic Optimization
 - Branch and Bound ...

Newton-Raphson Method

Iterative procedure to find a **root**

Newton-Raphson Method

Iterative procedure to find a **root**

Often solving for x when $f(x) = 0$ is hard \rightsquigarrow complicated function

Newton-Raphson Method

Iterative procedure to find a **root**

Often solving for x when $f(x) = 0$ is hard \rightsquigarrow complicated function

Solving for x when $f(x)$ is linear \rightsquigarrow easy

Newton-Raphson Method

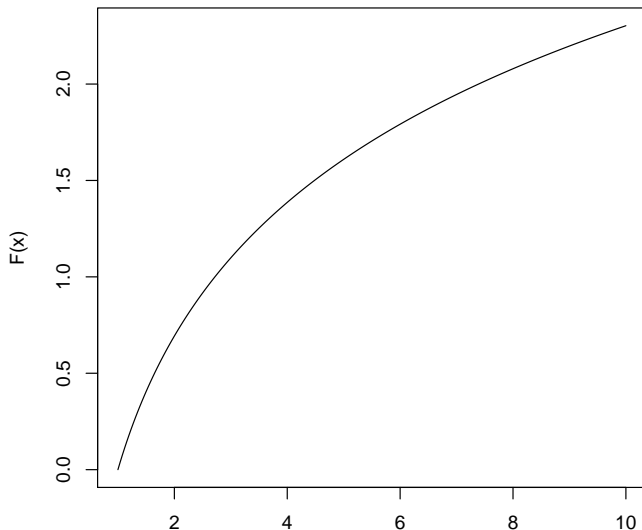
Iterative procedure to find a **root**

Often solving for x when $f(x) = 0$ is hard \rightsquigarrow complicated function

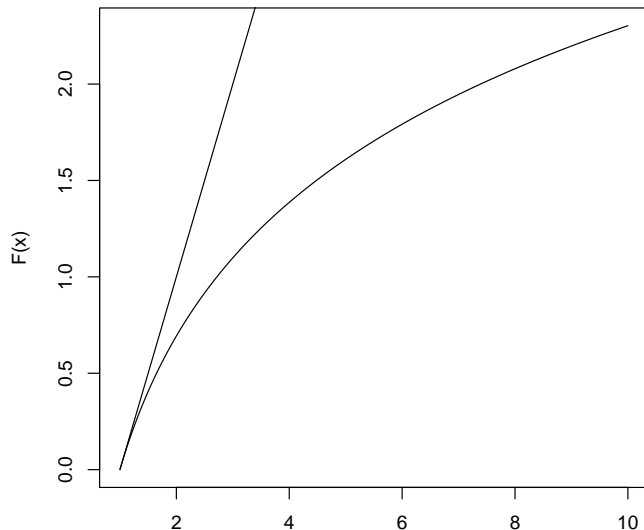
Solving for x when $f(x)$ is linear \rightsquigarrow easy

Approximate with **tangent line**, iteratively update

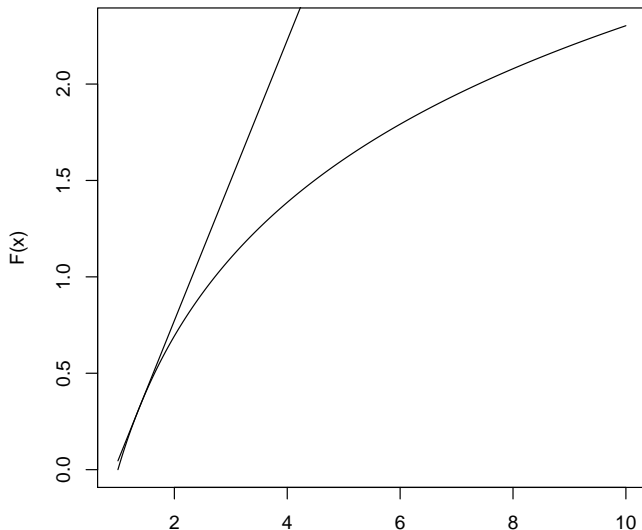
Tangent Line



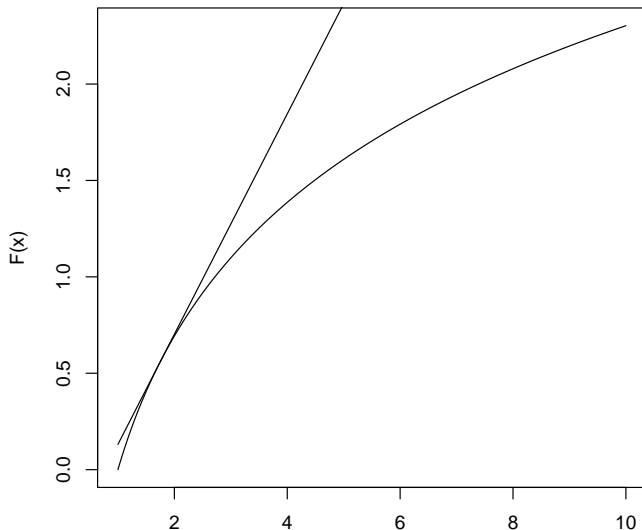
Tangent Line



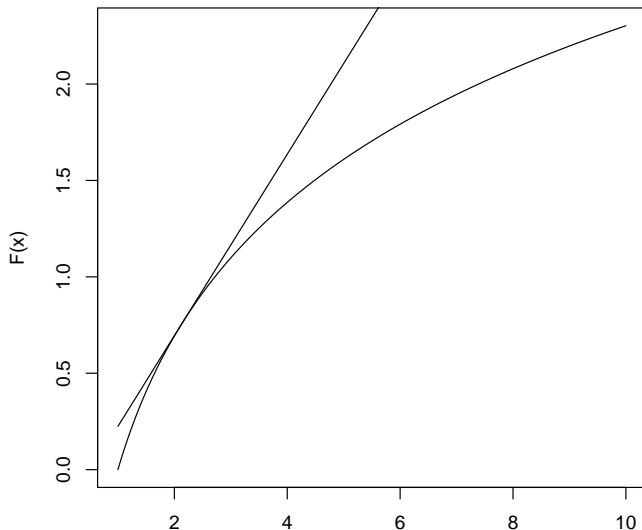
Tangent Line



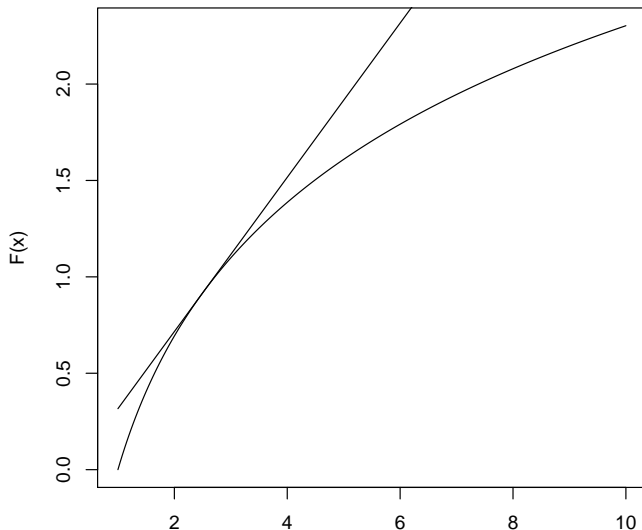
Tangent Line



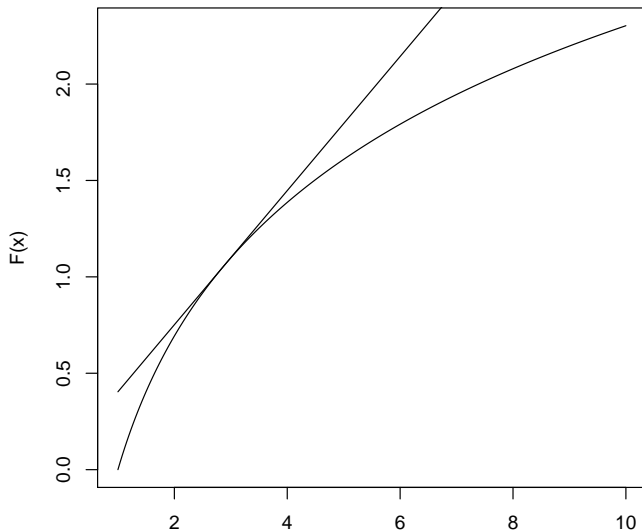
Tangent Line



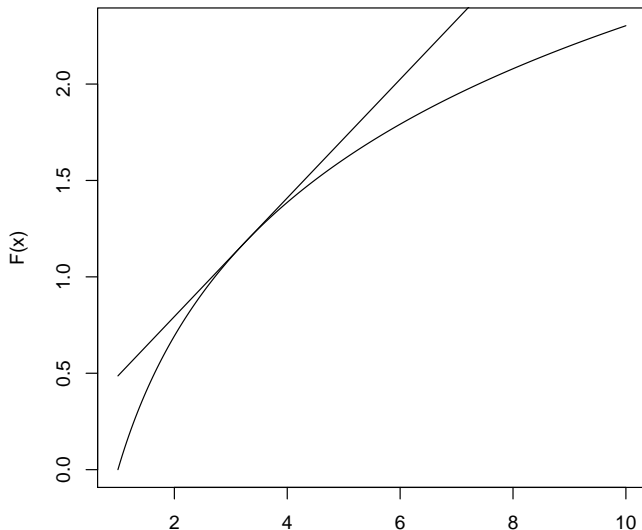
Tangent Line



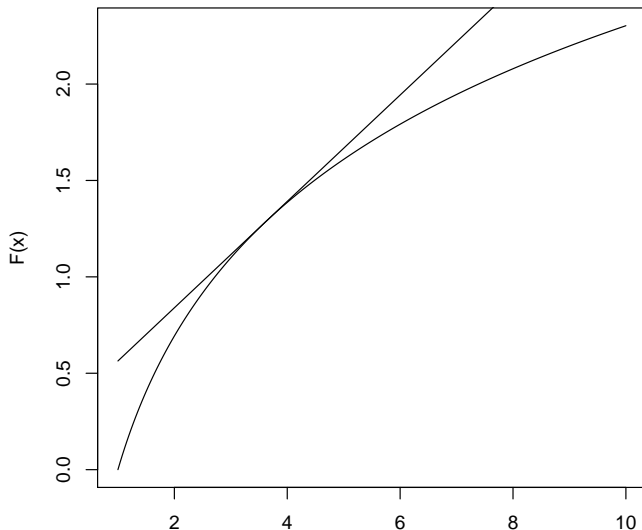
Tangent Line



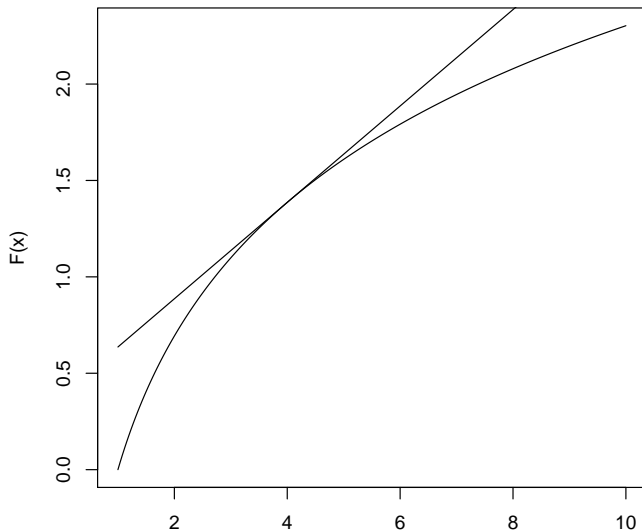
Tangent Line



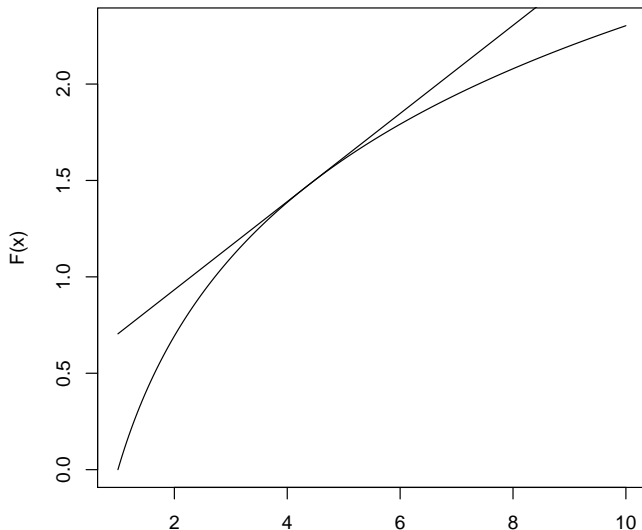
Tangent Line



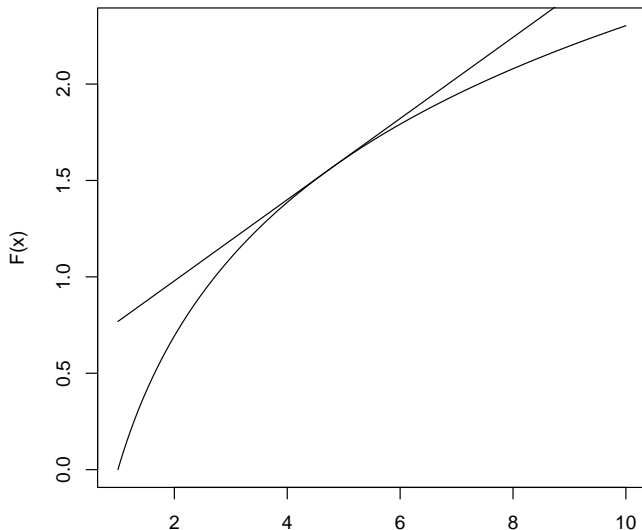
Tangent Line



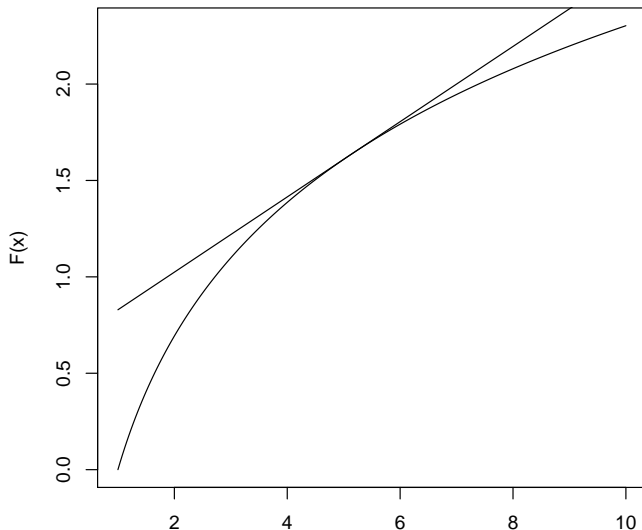
Tangent Line



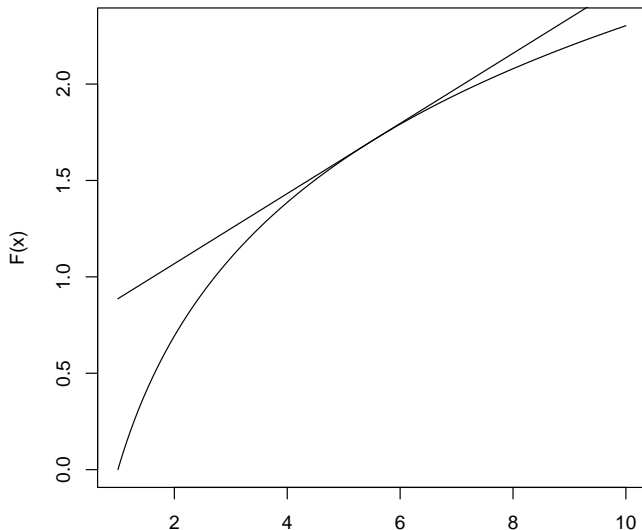
Tangent Line



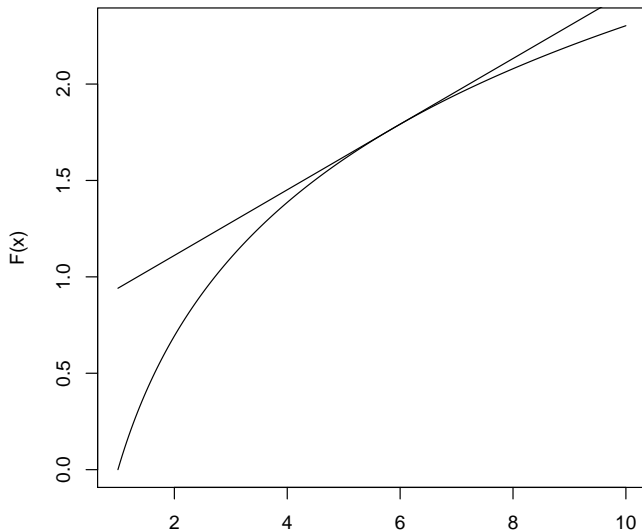
Tangent Line



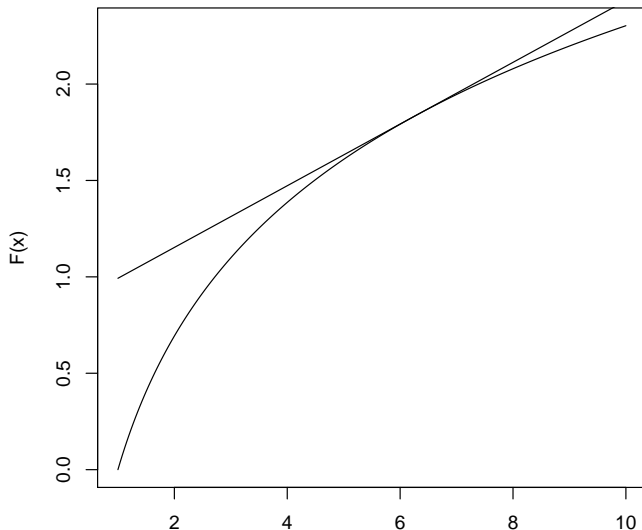
Tangent Line



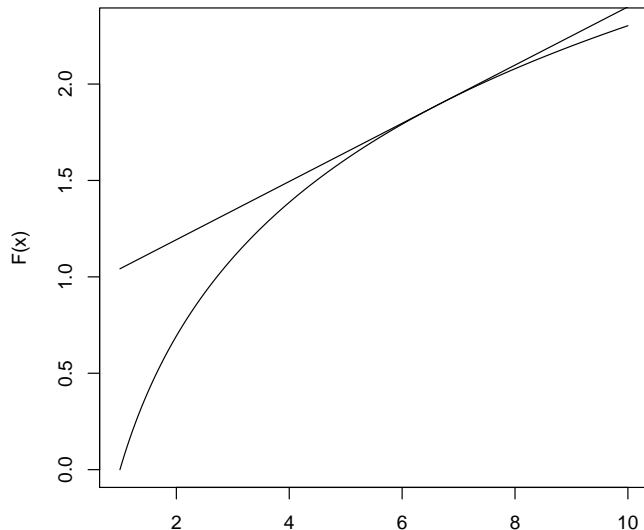
Tangent Line



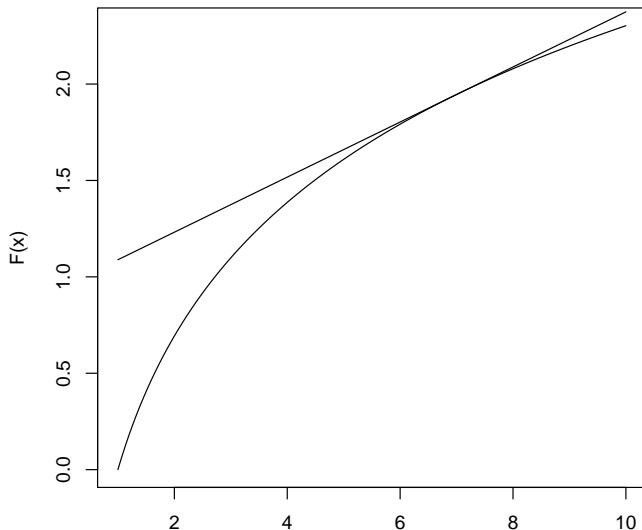
Tangent Line



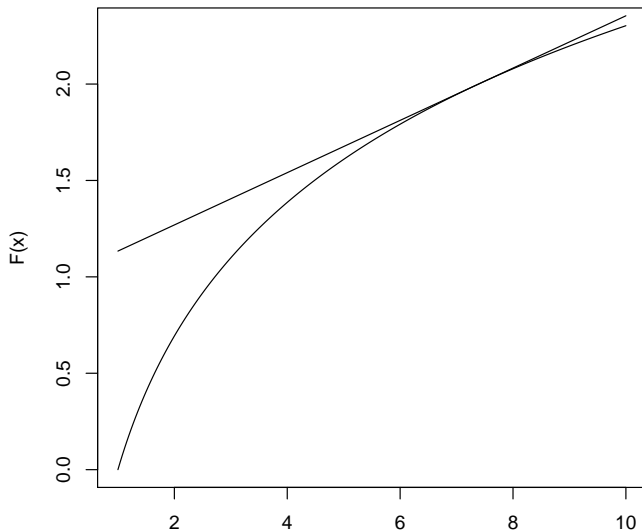
Tangent Line



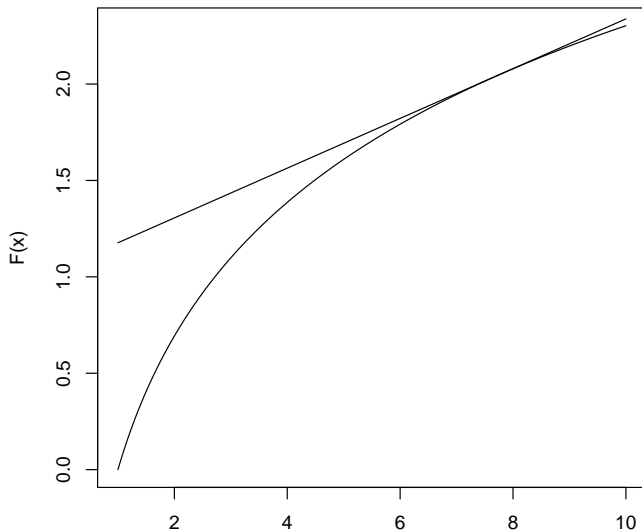
Tangent Line



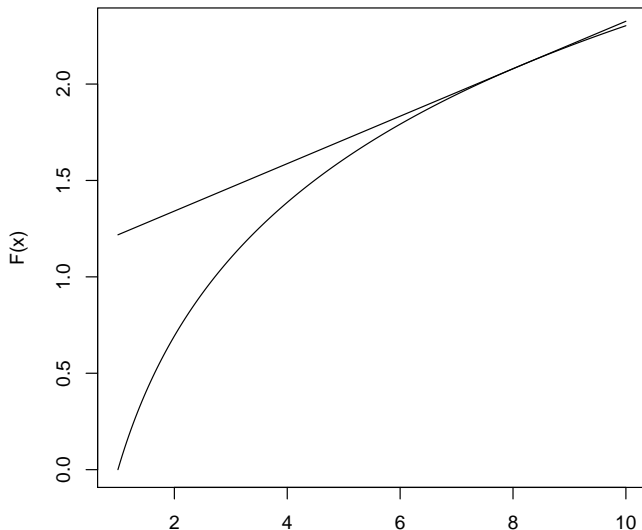
Tangent Line



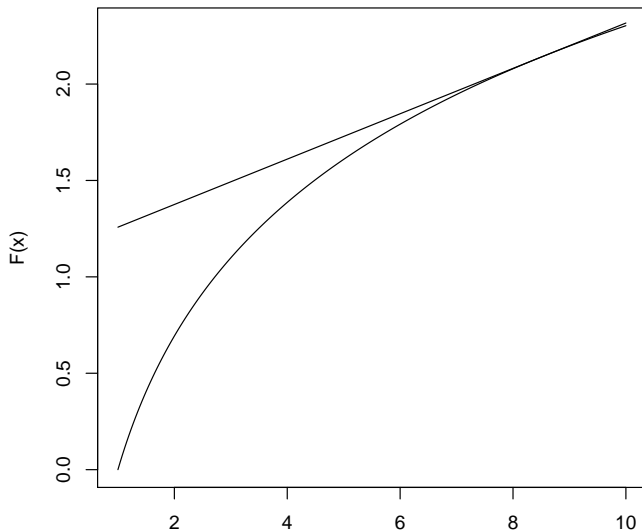
Tangent Line



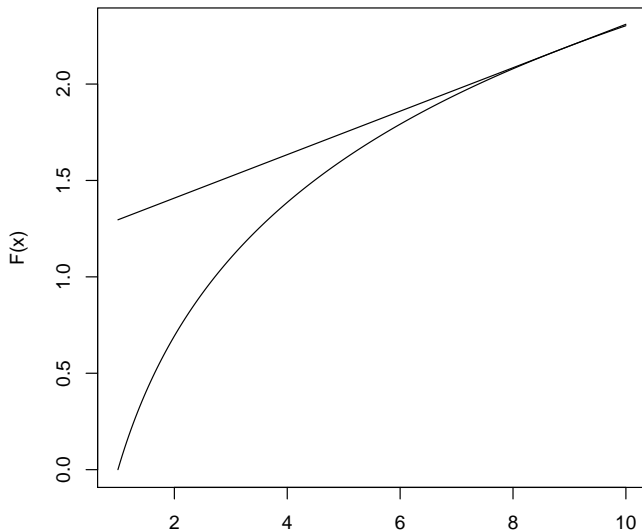
Tangent Line



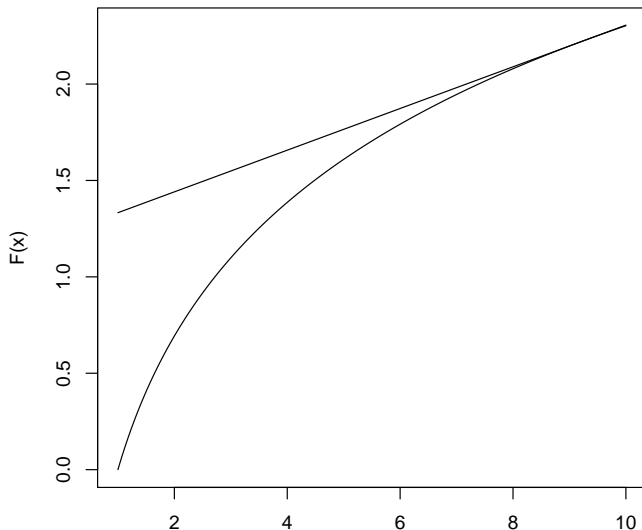
Tangent Line



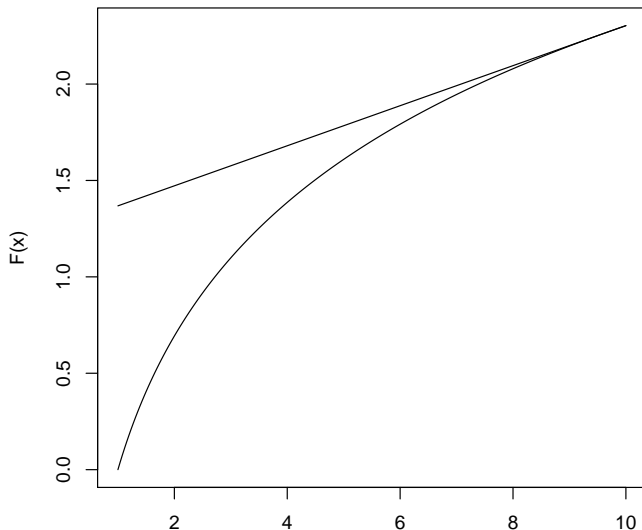
Tangent Line



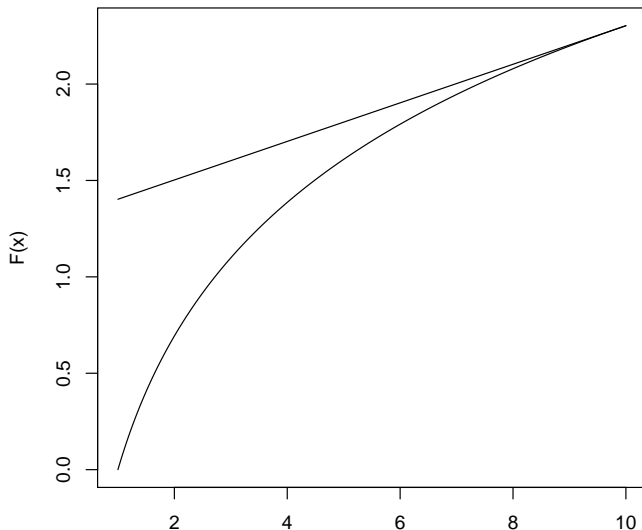
Tangent Line



Tangent Line



Tangent Line



Tangent Line

Formula for Tangent line at x_0 :

Tangent Line

Formula for Tangent line at x_0 :

$$g(x) = f'(x_0)(x - x_0) + f(x_0)$$

Tangent Line

Formula for Tangent line at x_0 :

$$g(x) = f'(x_0)(x - x_0) + f(x_0)$$

Tangent Line

Formula for Tangent line at x_0 :

$$g(x) = f'(x_0)(x - x_0) + f(x_0)$$

Tangent Line

Formula for Tangent line at x_0 :

$$g(x) = f'(x_0)(x - x_0) + f(x_0)$$

Newton-Raphson Method

Suppose we have some initial guess x_0 . We're going to approximate $f'(x)$ with the tangent line to generate a new guess

Newton-Raphson Method

Suppose we have some initial guess x_0 . We're going to approximate $f'(x)$ with the tangent line to generate a new guess

$$g(x) = f''(x_0)(x - x_0) + f'(x_0)$$

Newton-Raphson Method

Suppose we have some initial guess x_0 . We're going to approximate $f'(x)$ with the tangent line to generate a new guess

$$\begin{aligned}g(x) &= f''(x_0)(x - x_0) + f'(x_0) \\ 0 &= f''(x_0)(x_1 - x_0) + f'(x_0)\end{aligned}$$

Newton-Raphson Method

Suppose we have some initial guess x_0 . We're going to approximate $f'(x)$ with the tangent line to generate a new guess

$$\begin{aligned}g(x) &= f''(x_0)(x - x_0) + f'(x_0) \\0 &= f''(x_0)(x_1 - x_0) + f'(x_0) \\x_1 &= x_0 - \frac{f'(x_0)}{f''(x_0)}\end{aligned}$$

Newton-Raphson Method

Suppose we have some initial guess x_0 . We're going to approximate $f'(x)$ with the tangent line to generate a new guess

$$\begin{aligned}g(x) &= f''(x_0)(x - x_0) + f'(x_0) \\0 &= f''(x_0)(x_1 - x_0) + f'(x_0) \\x_1 &= x_0 - \frac{f'(x_0)}{f''(x_0)} \\x_{t+1} &= x_t - \frac{f'(x_t)}{f''(x_t)}\end{aligned}$$

Newton-Raphson Method

Suppose we have some initial guess x_0 . We're going to approximate $f'(x)$ with the tangent line to generate a new guess

$$\begin{aligned}g(x) &= f''(x_0)(x - x_0) + f'(x_0) \\0 &= f''(x_0)(x_1 - x_0) + f'(x_0) \\x_1 &= x_0 - \frac{f'(x_0)}{f''(x_0)} \\x_{t+1} &= x_t - \frac{f'(x_t)}{f''(x_t)}\end{aligned}$$

Perform iteratively until change in $|x_{t+1} - x_t| < \text{threshold}$

Newton-Raphson Method

Suppose we have some initial guess x_0 . We're going to approximate $f'(x)$ with the tangent line to generate a new guess

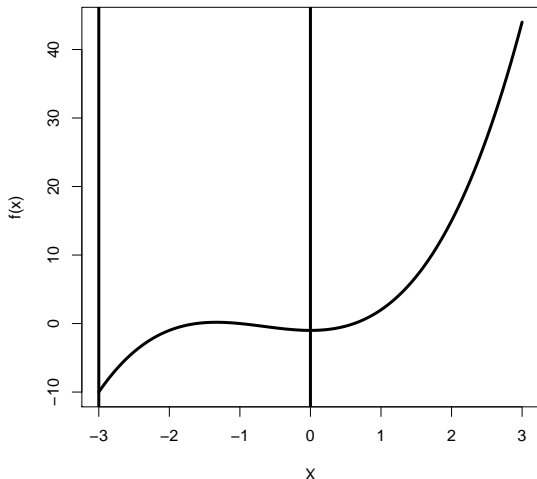
$$\begin{aligned}g(x) &= f''(x_0)(x - x_0) + f'(x_0) \\0 &= f''(x_0)(x_1 - x_0) + f'(x_0) \\x_1 &= x_0 - \frac{f'(x_0)}{f''(x_0)} \\x_{t+1} &= x_t - \frac{f'(x_t)}{f''(x_t)}\end{aligned}$$

Perform iteratively until change in $|x_{t+1} - x_t| < \text{threshold}$

Example Function

$f(x) = x^3 + 2x^2 - 1$ find x that maximizes $f(x)$ with $x \in [-3, 0]$

$$x^3 + 2x^2 - 1$$



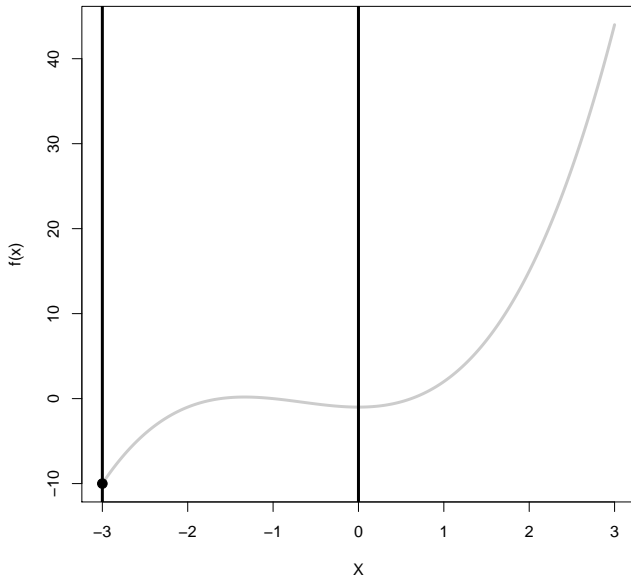
$$f'(x) = 3x^2 + 4x$$

$$f''(x) = 6x + 4$$

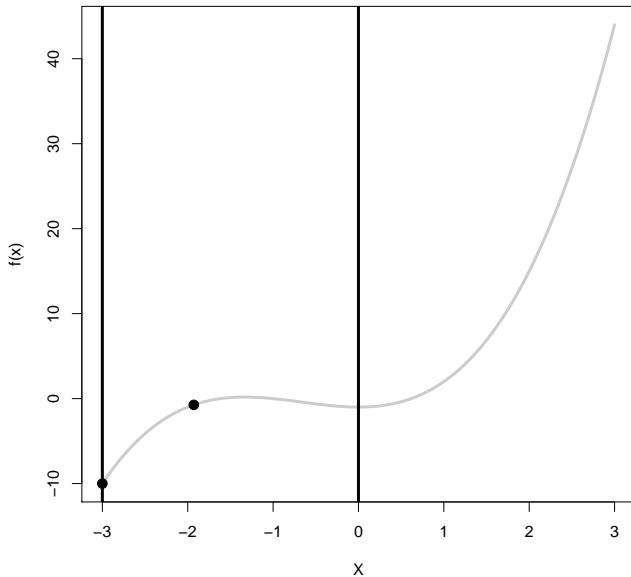
Suppose we have guess x_t then the next step is:

$$x_{t+1} = x_t - \frac{3x_t^2 + 4x_t}{6x_t + 4}$$

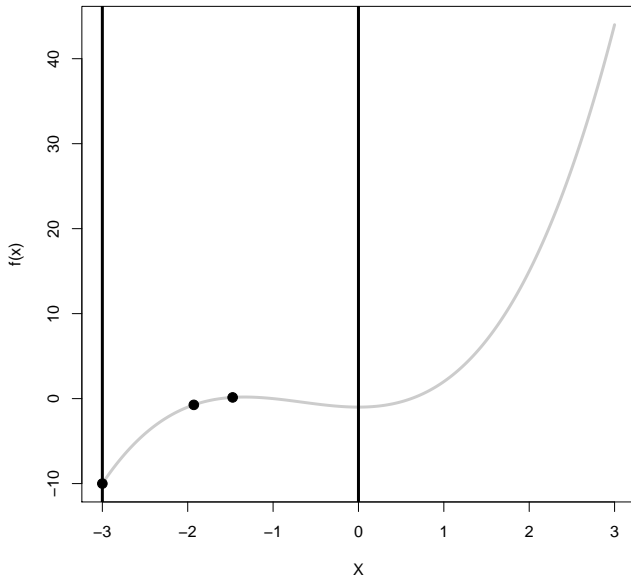
$$x^3 + 2x^2 - 1$$



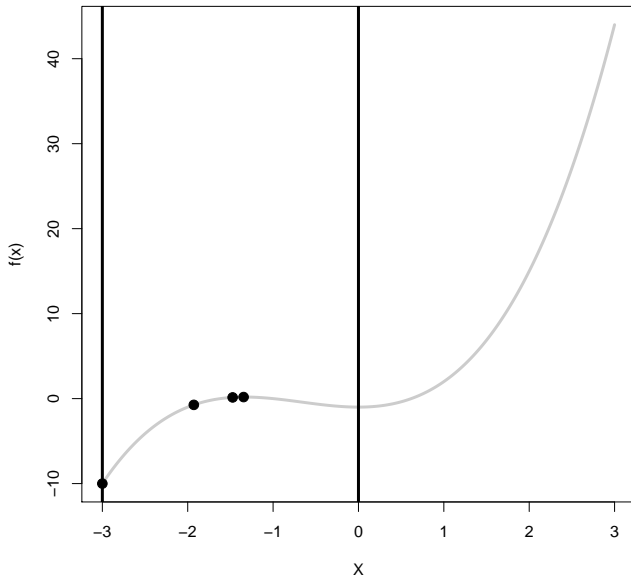
$$x^3 + 2x^2 - 1$$



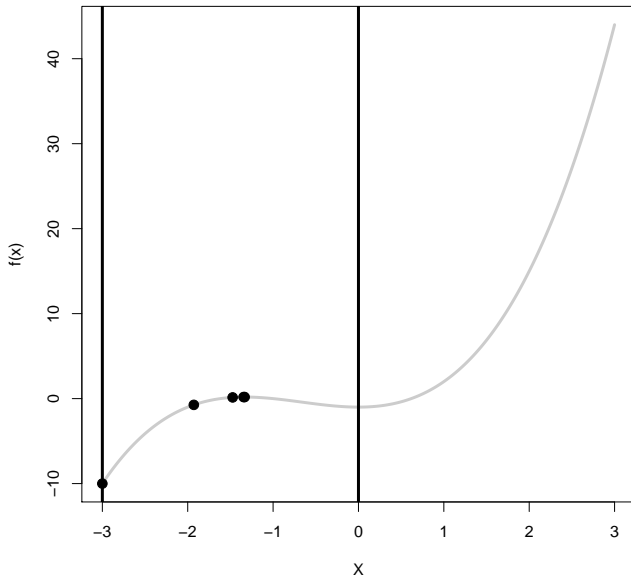
$$x^3 + 2x^2 - 1$$



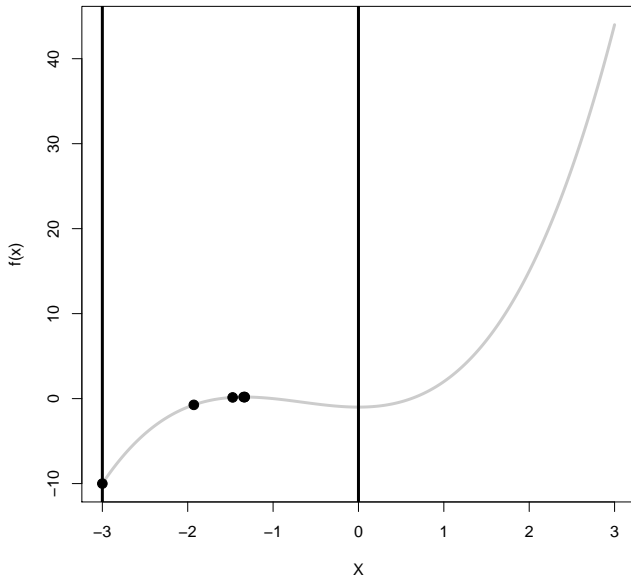
$$x^3 + 2x^2 - 1$$



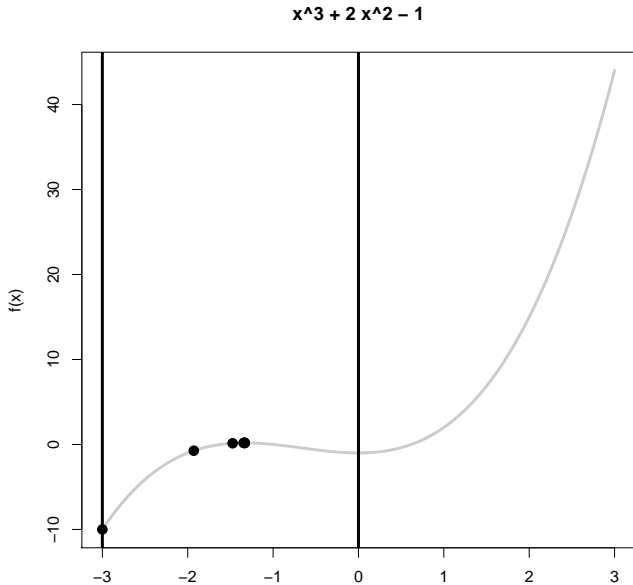
$$x^3 + 2x^2 - 1$$



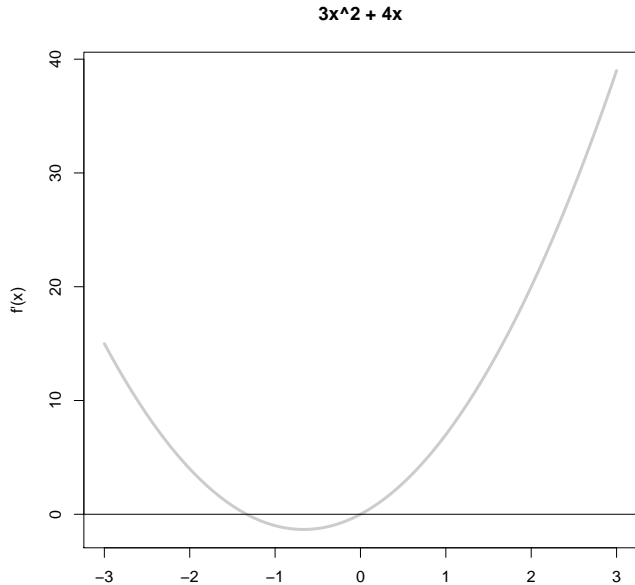
$$x^3 + 2x^2 - 1$$



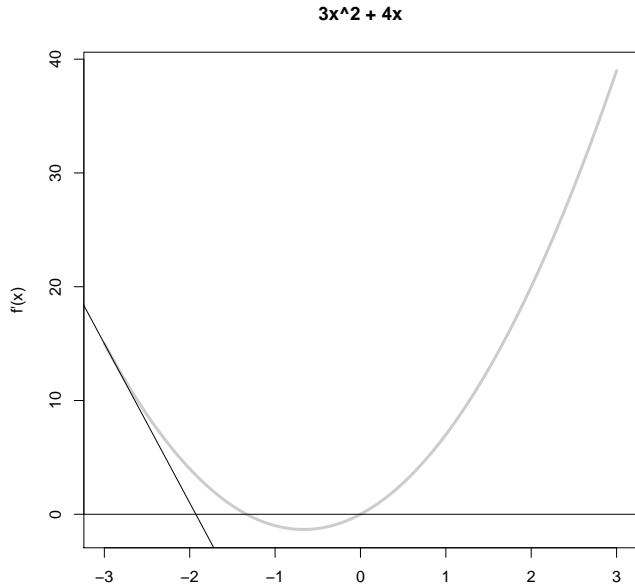
$$x^* = -1.3333$$



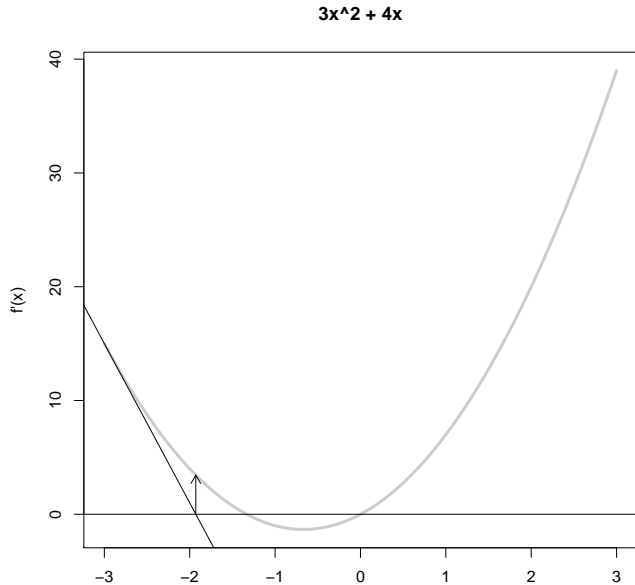
What is Happening with the Roots



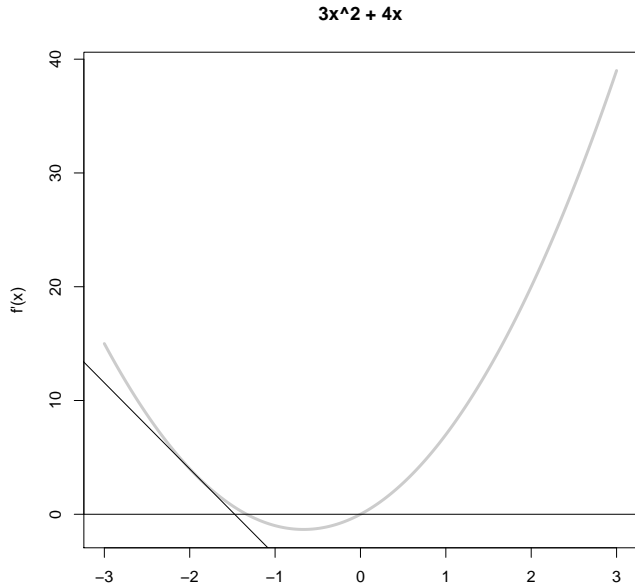
What is Happening with the Roots



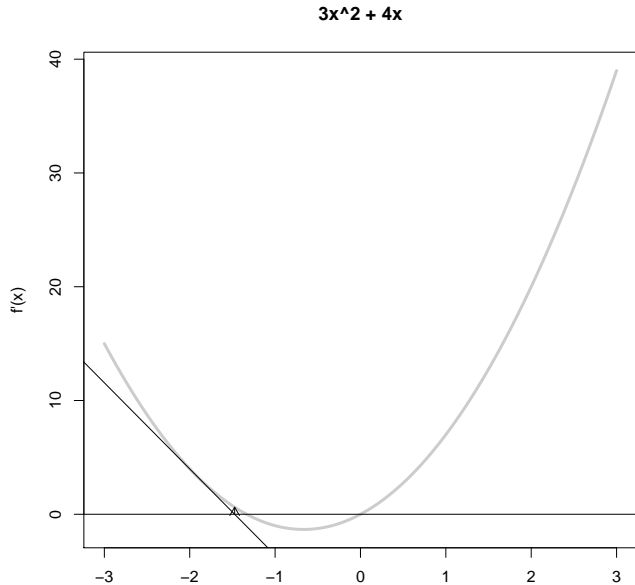
What is Happening with the Roots



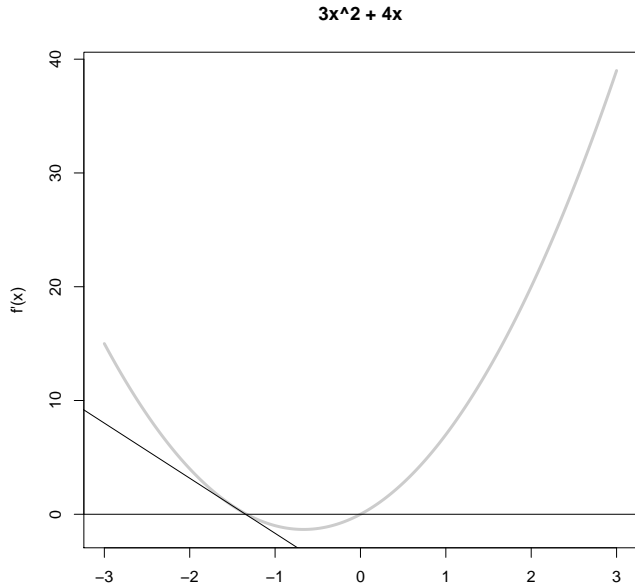
What is Happening with the Roots



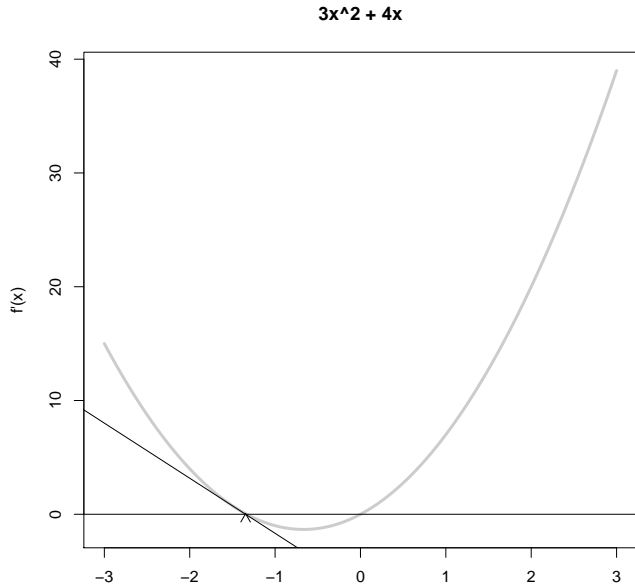
What is Happening with the Roots



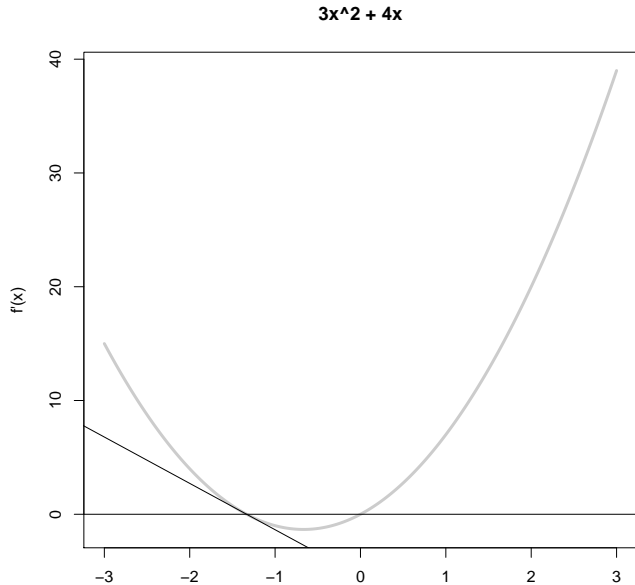
What is Happening with the Roots



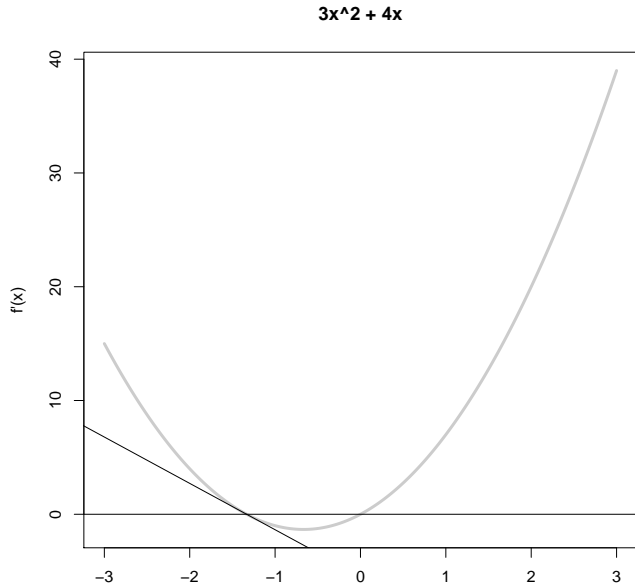
What is Happening with the Roots



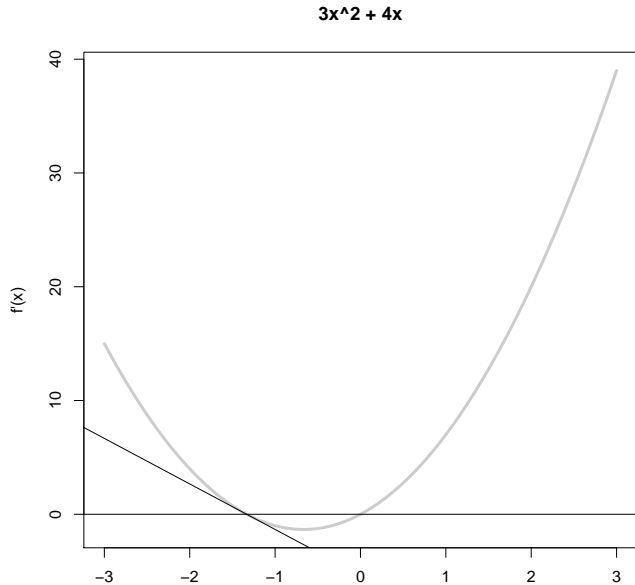
What is Happening with the Roots



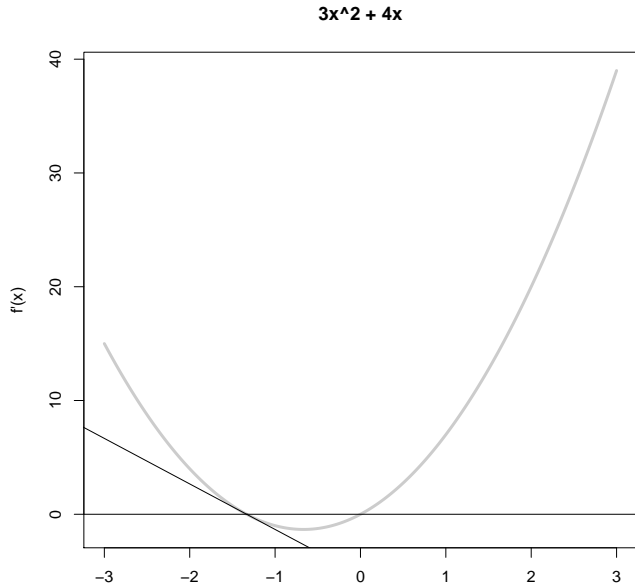
What is Happening with the Roots



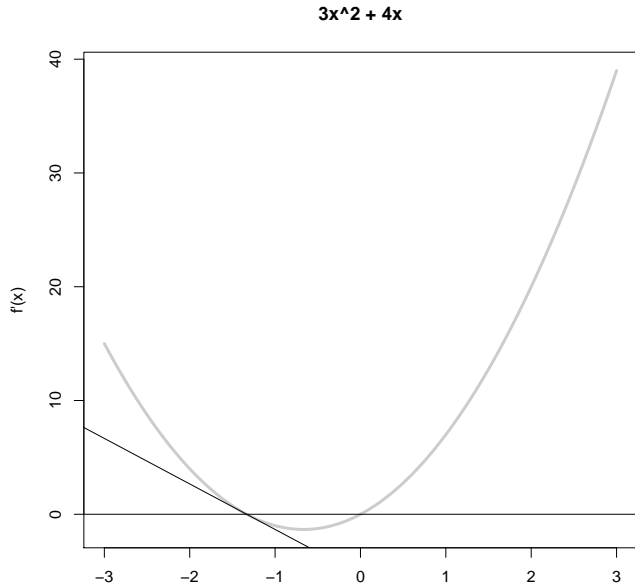
What is Happening with the Roots



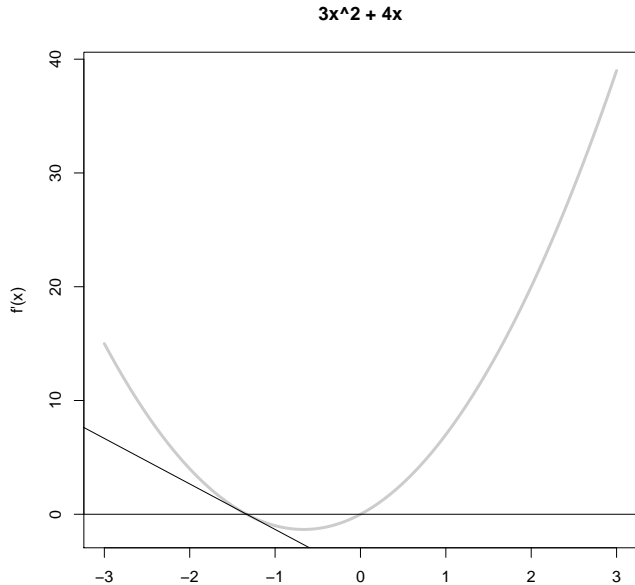
What is Happening with the Roots



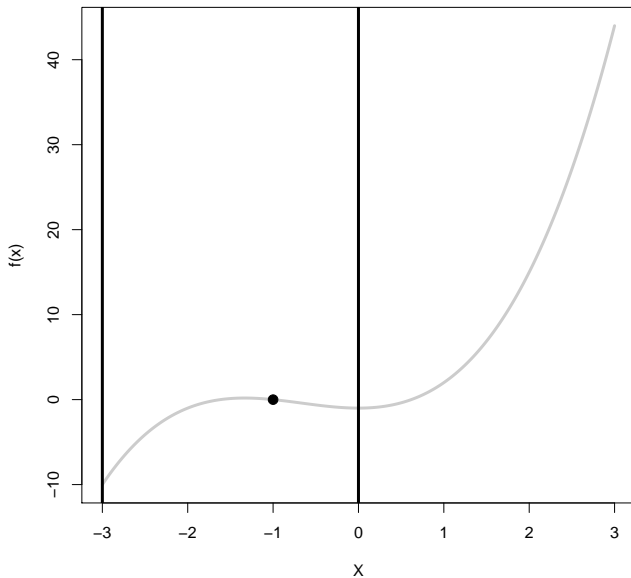
What is Happening with the Roots



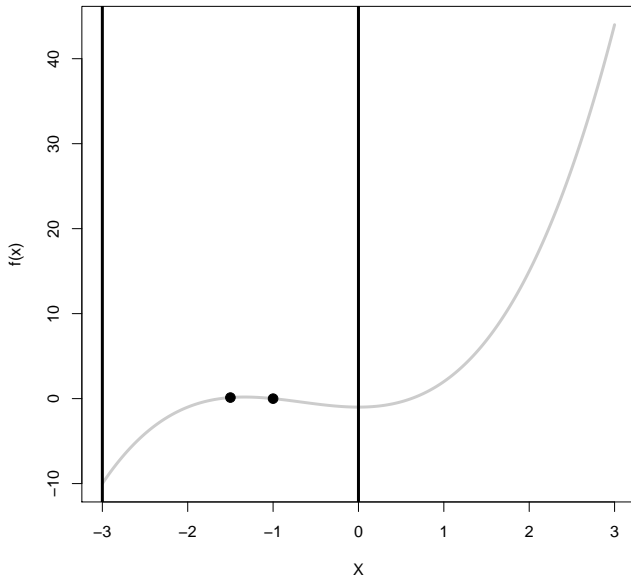
What is Happening with the Roots



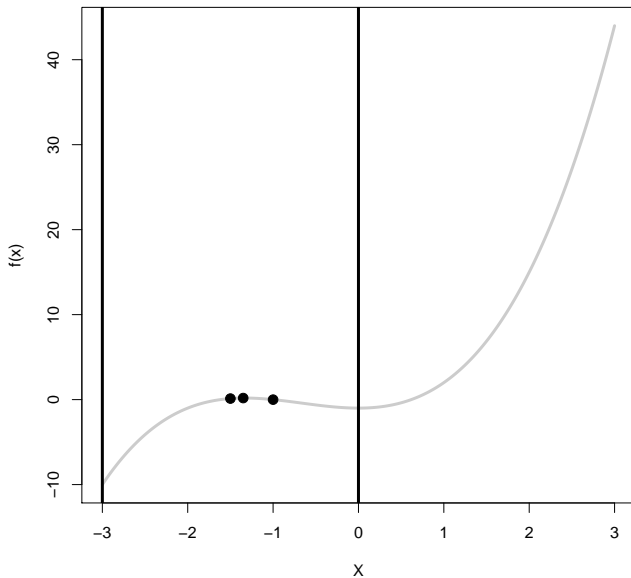
$$x^3 + 2x^2 - 1$$



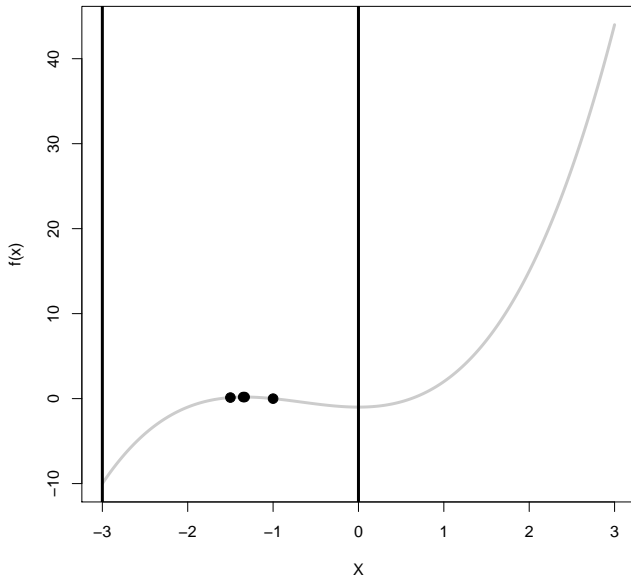
$$x^3 + 2x^2 - 1$$



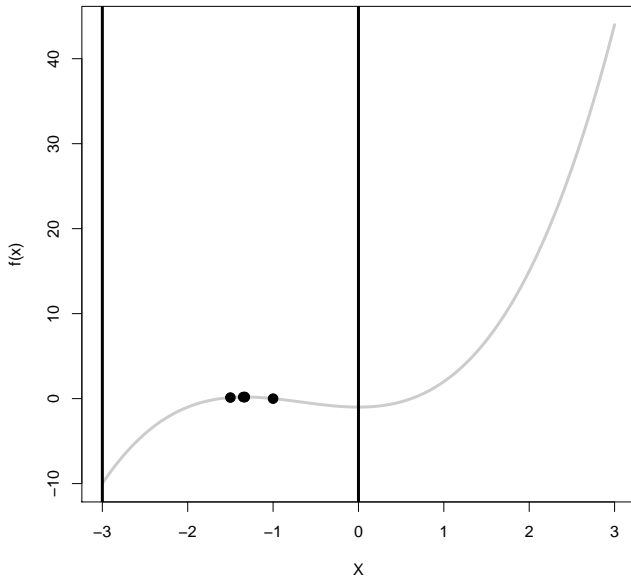
$$x^3 + 2x^2 - 1$$



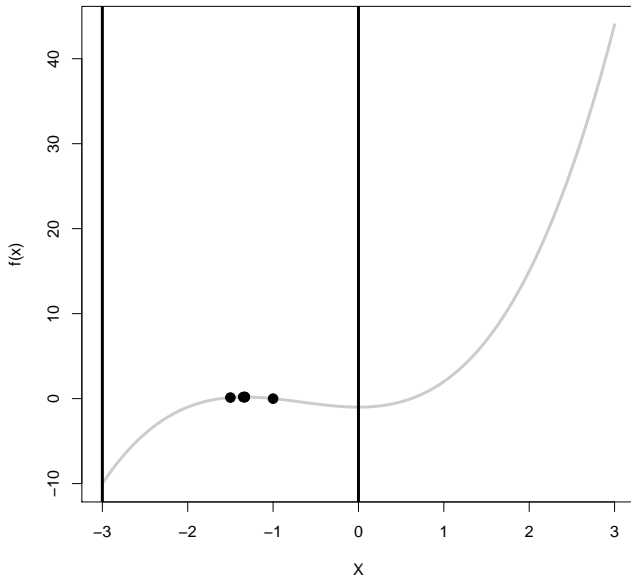
$$x^3 + 2x^2 - 1$$



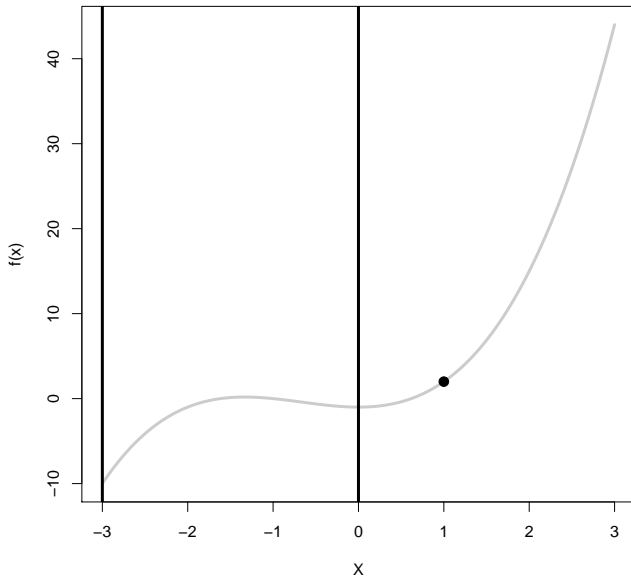
$$x^3 + 2x^2 - 1$$



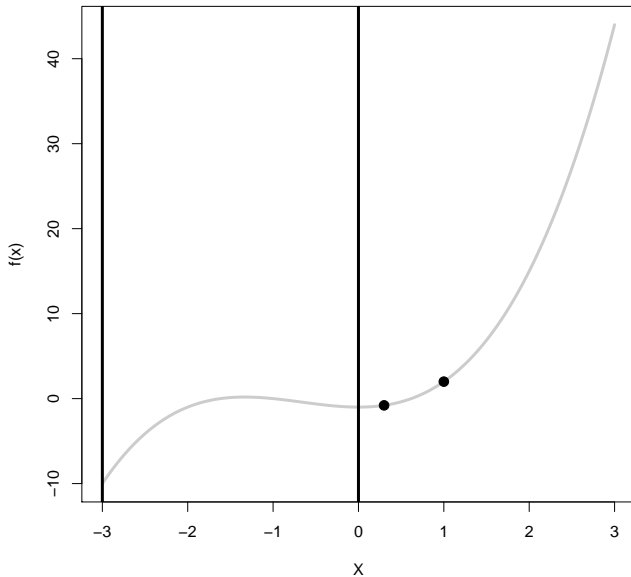
$$x^3 + 2x^2 - 1$$



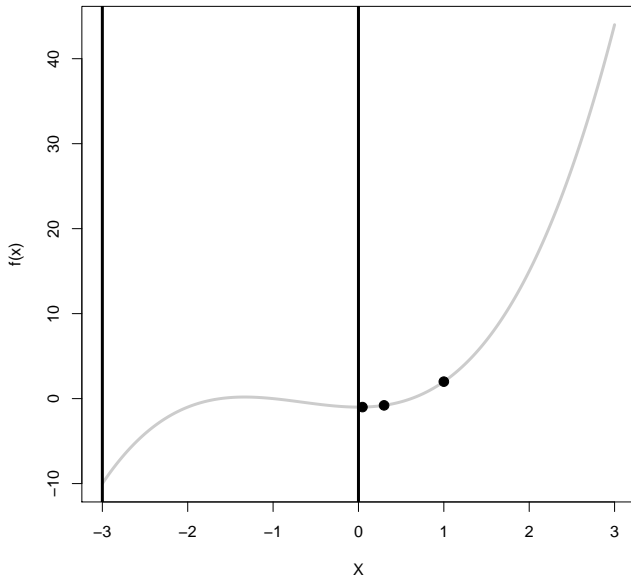
$$x^3 + 2x^2 - 1$$



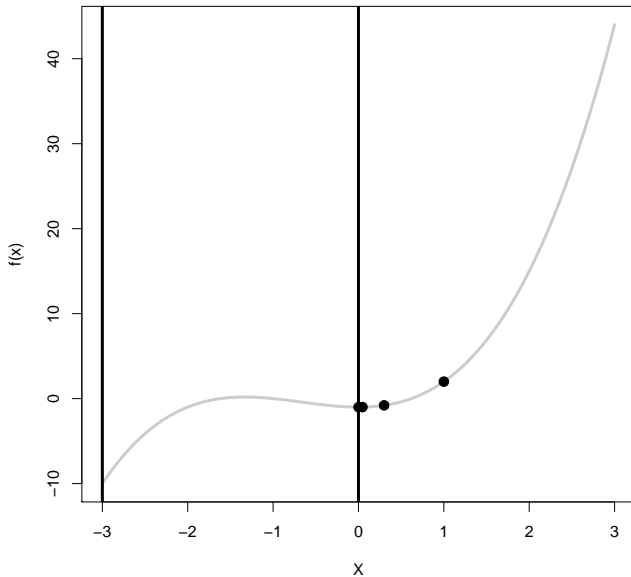
$$x^3 + 2x^2 - 1$$



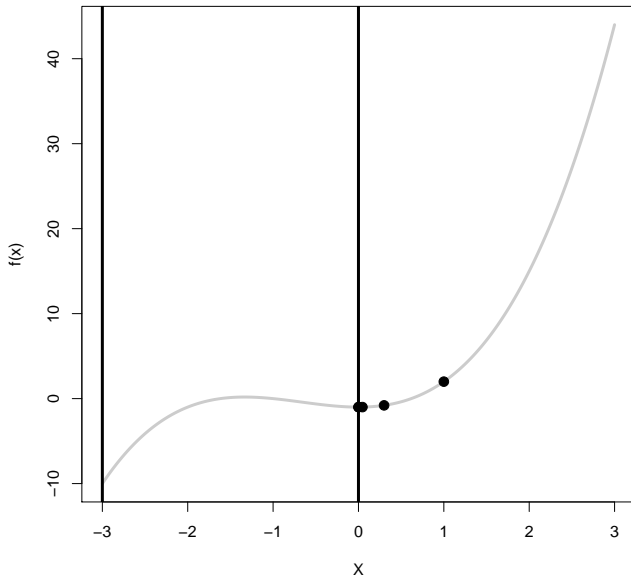
$$x^3 + 2x^2 - 1$$



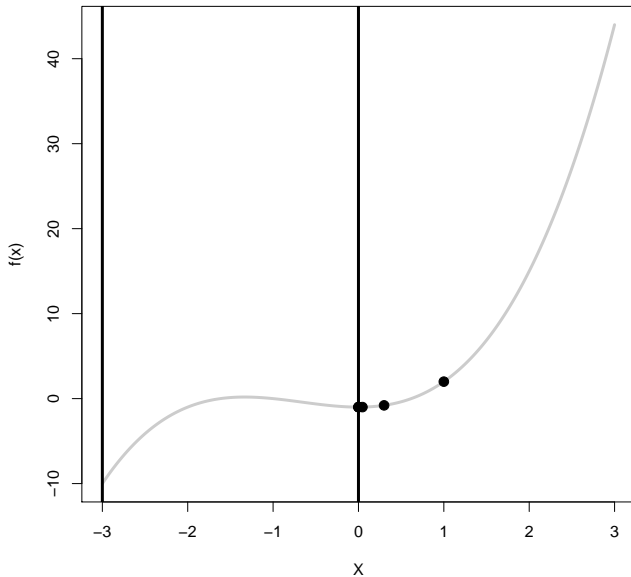
$$x^3 + 2x^2 - 1$$



$$x^3 + 2x^2 - 1$$



$$x^3 + 2x^2 - 1$$



Multivariate Optimization

$$\log(p(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X})) = \sum_{i=1}^N [y_i \log(\Phi(\mathbf{X}_i\boldsymbol{\beta})) + (1 - y_i) \log(1 - \Phi(\mathbf{X}_i\boldsymbol{\beta}))] + c$$

To do so:

Apply **BFGS** (quasi-Newton) in R, in `optim`

R code

Multivariate Optimization

$$\log(p(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X})) = \sum_{i=1}^N [y_i \log(\Phi(\mathbf{X}_i\boldsymbol{\beta})) + (1 - y_i) \log(1 - \Phi(\mathbf{X}_i\boldsymbol{\beta}))] + c$$

To do so:

Apply **BFGS** (quasi-Newton) in R, in `optim`

R code

Estimates: predict, classify, describe, ...

Probit Regression, with Prior

Consider the following data generation process

$$Y_i \sim \text{Bernoulli}(\pi_i)$$

$$\pi_i = \Phi(\mathbf{X}_i\boldsymbol{\beta})$$

$$\beta_j \sim \text{Normal}(\mu, \sigma^2)$$

$$\begin{aligned} p(\boldsymbol{\beta}|\mathbf{X}, \mathbf{Y}, \mu, \sigma^2) &\propto p(\boldsymbol{\beta}|\mu, \sigma^2) \times p(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}) \\ &\propto \prod_{j=1}^2 \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\beta_j - \mu)^2}{2\sigma^2}\right) \times \prod_{i=1}^N \Phi(\mathbf{X}_i\boldsymbol{\beta})^{y_i} (1 - \Phi(\mathbf{X}_i\boldsymbol{\beta})) \end{aligned}$$

Homework \rightsquigarrow explore behavior of $\hat{\boldsymbol{\beta}}$ as μ, σ^2 vary.

Where We're Going

- 1) Task
- 2) Objective Function
- 3) Optimization procedure

All supposes we have data.

Next week \rightsquigarrow converting text data